

Intermediate College Physics

Richard Fitzpatrick

Professor of Physics
The University of Texas at Austin

Contents

Table of Contents	1
Preface	3
1 Newtonian Dynamics	5
1.1 Units	5
1.1.1 Mks Units	5
1.1.2 Standard Prefixes	6
1.1.3 Other Units	7
1.1.4 Dimensional Analysis	8
1.1.5 Deriving Physical Relationships	8
1.1.6 Scaling Laws	9
1.2 Newton's Laws of Motion	10
1.2.1 Introduction	10
1.2.2 Newton's First Law of Motion	10
1.2.3 Newton's Second Law of Motion	11
1.2.4 Newton's Third Law of Motion	12
1.3 Motion of Single Particle	13
1.3.1 Impulse	13
1.3.2 Work	14
1.3.3 Conservative Forces	15
1.3.4 Potential Energy	16
1.3.5 Energy Conservation	17
1.3.6 Energy Diagrams	17
1.4 Motion of System of Many Particles	20
1.4.1 Equations of Motion	20
1.4.2 Center of Mass	21

1.4.3	Explosion of Krypton	22
1.4.4	Conservation of Linear Momentum	22
1.4.5	Conservation of Angular Momentum	23
1.5	Invariance Laws	24
1.5.1	Inertial Reference Frame	24
1.5.2	Rotational Invariance	25
1.5.3	Translational Invariance	25
1.5.4	Galilean Invariance	26
1.6	Two-Particle Collisions	28
1.6.1	One-Dimensional Collisions	28
1.6.2	Totally Inelastic Collisions	32
1.6.3	Two-Dimensional Collisions	32
1.7	Rigid Body Rotation	36
1.7.1	Fundamental Equations	36
1.7.2	Moment of Inertia Tensor	37
1.7.3	Rotational Kinetic Energy	39
1.7.4	Power	39
1.7.5	Uniform Flywheel	40
1.7.6	Gravitational Collapse of Star	42
1.7.7	Gyroscopic Precession	43
1.8	Newtonian Gravity	45
1.8.1	Gravity	45
1.8.2	Gauss's Law	45
1.8.3	Gravitational Field of Earth	47
1.8.4	Gravitational Potential Energy	48
1.8.5	Gravitational Potential	50
1.9	Planetary Motion	51
1.9.1	Kepler's Laws	51
1.9.2	Planetary Conservation Laws	51
1.9.3	Plane Polar Coordinates	53
1.9.4	Conic Sections	54
1.9.5	Kepler's Second Law	57
1.9.6	Kepler's First Law	58
1.9.7	Kepler's Third Law	59
1.9.8	Orbital Energies	60
1.10	Spheroidal Mass Distributions	61
1.10.1	Gravitational Potential of Uniform Spheroid	61
1.10.2	Rotational Flattening of Earth	64
1.10.3	Surface Gravity of Earth	67
1.10.4	MacCullagh's Formula	68
1.10.5	Gravitational Torque on Axisymmetric Mass Distribution	69
1.10.6	Lunisolar Precession of Earth	70

1.10.7	Two-Body Dynamics	76
1.10.8	Binary Star Systems	78
1.10.9	Tidal Elongation of Earth	80
2	Classical Electromagnetism	85
2.1	Electrostatic Fields	85
2.1.1	Electricity	85
2.1.2	Coulomb's Law	86
2.1.3	Electric Field	88
2.1.4	Electric Scalar Potential	90
2.1.5	Electric Potential Energy	91
2.1.6	Gauss's Law	92
2.1.7	Applications of Gauss's Law	97
2.1.8	Electrostatic Energy	99
2.1.9	Poisson's Equation	104
2.1.10	Uniqueness Theorem	104
2.1.11	Ohm's Law	105
2.1.12	Ideal Conductors	107
2.1.13	Capacitors	111
2.1.14	Method of Images	115
2.2	Magnetostatic Fields	123
2.2.1	Magnetism	123
2.2.2	Magnetic Field	123
2.2.3	Ampère's Law	126
2.2.4	Lorentz Force Law	127
2.2.5	Charged Particle Motion in a Magnetic Field	130
2.2.6	Hall Effect	132
2.2.7	Biot-Savart Law	133
2.2.8	Magnetic Vector Potential	136
2.2.9	Magnetic Monopoles	138
2.2.10	Ampère's Circuital Law	139
2.2.11	Magnetic Field of a Solenoid	142
2.3	Magnetic Induction	143
2.3.1	Faraday's Law	143
2.3.2	Electric Scalar Potential	147
2.3.3	Gauge Invariance	147
2.3.4	Inductance	149
2.3.5	Self Inductance	151
2.3.6	<i>RC</i> Circuits	153
2.3.7	Mutual Inductance	156
2.3.8	Magnetic Energy	158
2.3.9	Motional Emf	163
2.3.10	Alternating Current Generators	165

2.3.11	Alternating Current Circuits	169
2.3.12	Alternating Current Motors	172
2.3.13	Transformers	174
2.4	Maxwell's Equations	178
2.4.1	Displacement Current	178
2.4.2	Maxwell's Equations	183
2.4.3	Potential Formulation of Maxwell's Equations	184
2.4.4	Electromagnetic Waves	185
2.4.5	Energy Conservation	191
2.4.6	Electromagnetic Momentum	194
3	Special Relativity	199
3.1	Experimental Basis of Special Relativity	199
3.1.1	Sound Waves in a Gas	199
3.1.2	Light Waves in a Vacuum	201
3.1.3	Aberration of Starlight and Stellar Parallax	202
3.1.4	Fizeau and Airy Experiments	208
3.1.5	Michelson-Morley Experiment	208
3.1.6	Lorentz-Fitzgerald Contraction	209
3.1.7	Kennedy-Thorndike Experiment	209
3.2	Theoretical Basis of Special Relativity	210
3.2.1	Postulates of Special Relativity	210
3.2.2	Invariance of Transverse Lengths	212
3.2.3	Time Dilation	212
3.2.4	Length Contraction	214
3.2.5	Clock Error	216
3.2.6	Galilean Transformation	217
3.2.7	Lorentz Transformation	218
3.2.8	Spacetime Interval	219
3.2.9	Transformation of Velocity	221
3.2.10	Causality	222
3.2.11	Relativistic Aberration of Light	223
3.2.12	Relativistic Beaming of Light	224
3.2.13	Light Propagation through Dielectric Media	225
3.3	Relativistic Dynamics	225
3.3.1	Transformation of Acceleration	225
3.3.2	Relativistic Equation of Motion	227
3.3.3	Work and Energy	228
3.3.4	Relativistic Energy	229
3.3.5	Relativistic Energy-Momentum Relation	230
3.3.6	Transformation of Energy and Momentum	230
3.3.7	Relativistic Momentum Conservation	231
3.3.8	Photons	232

3.3.9	Relativistic Doppler Effect	233
3.3.10	Transverse Doppler Effect	234
3.3.11	Compton Scattering	236
3.3.12	Relativistic Inelastic Scattering	238
3.4	Relativity and Electromagnetism	239
3.4.1	Transformation of Electromagnetic Fields	239
3.4.2	Electromagnetic Fields of a Moving Charge	245
4	Quantum Mechanics	249
4.1	Experimental Basis of Quantum Mechanics	249
4.1.1	Wave-Particle Duality	249
4.1.2	Photoelectric Effect	250
4.1.3	Compton Scattering	251
4.1.4	Photon Polarization	252
4.1.5	Double-Slit Interference of Light	253
4.1.6	Electron Diffraction	254
4.1.7	Helium Diffraction	254
4.1.8	Two-Source Particle Interference	254
4.1.9	de Broglie's Hypothesis	255
4.2	Wave Mechanics	255
4.2.1	Wavefunctions	255
4.2.2	Schrödinger's Equation	256
4.2.3	Probability Interpretation of Wavefunction	257
4.2.4	Wave Packets	259
4.2.5	Group Velocity	262
4.2.6	Wave Dispersion	263
4.2.7	Heisenberg's Uncertainty Principle	264
4.2.8	Wavefunction Collapse	265
4.2.9	Stationary States	265
4.3	One-Dimensional Wave Mechanics	267
4.3.1	Particle in Infinite Square Potential Well	267
4.3.2	Particle in Finite Square Potential Well	270
4.3.3	Square Potential Barrier	274
4.3.4	WKB Approximation	278
4.3.5	Cold Emission	281
4.3.6	Alpha Decay	282
4.3.7	Simple Harmonic Oscillator	284
4.4	Three-Dimensional Wave Mechanics	287
4.4.1	Three-Dimensional Wave Mechanics	287
4.4.2	Particle in Box	288
4.4.3	Degenerate Electron Gas	290
4.4.4	White-Dwarf Star	292

5	Thermal Physics	295
5.1	Probability Theory	295
5.1.1	Probability	295
5.1.2	Binomial Probability Distribution	297
5.1.3	Mean, Variance, and Standard Deviation	298
5.1.4	Application to Binomial Probability Distribution	299
5.1.5	Random Walk	301
5.1.6	Continuous Probability Distribution	302
5.1.7	Gaussian Probability Distribution	303
5.2	Ideal Gas	307
5.2.1	Ideal Gas Law	307
5.2.2	First Law of Thermodynamics	308
5.2.3	Specific Heat Capacity	309
5.2.4	Isothermal and Adiabatic Expansion	310
5.2.5	Hydrostatic Equilibrium of Atmosphere	311
5.2.6	Isothermal Atmosphere	312
5.2.7	Adiabatic Atmosphere	313
5.2.8	Bulk Modulus	317
5.2.9	Sound Waves	318
5.3	Kinetic Theory	319
5.3.1	Fundamental Assumptions	319
5.3.2	Molecular Flux	319
5.3.3	Pressure	321
5.3.4	Law of Equipartition of Energy	322
5.3.5	Partial Pressure	322
5.3.6	Internal Energy	323
5.3.7	Brownian Motion	324
5.3.8	Mean Free Path	326
5.3.9	Diffusion	328
5.3.10	Thermal Conductivity	330
5.3.11	Viscosity	332
5.3.12	Molecular Flow	334
5.3.13	Molecular Effusion	335
5.4	Statistical Mechanics	336
5.4.1	Specification of State of Many-Particle System	336
5.4.2	Principle of Equal A Priori Probabilities	338
5.4.3	Probability Calculations	339
5.4.4	Number of Accessible States of Ideal Gas	340
5.4.5	Thermal Interaction	341
5.4.6	Thermodynamic Temperature	342
5.4.7	Boltzmann Probability Distribution	344
5.5	Applications of Statistical Mechanics	346

5.5.1	Two-State System	346
5.5.2	Spin-1/2 Paramagnetism	347
5.5.3	Adiabatic Demagnetization	350
5.5.4	Thermal Expansion	351
5.5.5	Equipartition Theorem	353
5.5.6	Harmonic Oscillator	355
5.5.7	Specific Heat Capacities	357
5.5.8	Specific Heats of Gases	358
5.5.9	Maxwell Velocity Distribution	362
5.6	Standing-Wave States	367
5.6.1	Counting Standing-Wave States	367
5.6.2	Planck Radiation Law	370
5.6.3	Black-Body Radiation	372
5.6.4	Stefan-Boltzmann Law	373
5.6.5	Specific Heats of Solids	375
5.6.6	Conduction Electrons in Metal	381
A	Vector Algebra and Vector Calculus	387
A.1	Introduction	387
A.2	Scalars and Vectors	387
A.3	Vector Algebra	388
A.4	Cartesian Components of a Vector	390
A.5	Coordinate Transformations	391
A.6	Scalar Product	392
A.7	Vector Area	394
A.8	Vector Product	395
A.9	Rotation	397
A.10	Scalar Triple Product	399
A.11	Vector Triple Product	400
A.12	Vector Calculus	401
A.13	Line Integrals	401
A.14	Vector Line Integrals	404
A.15	Surface Integrals	404
A.16	Vector Surface Integrals	406
A.17	Volume Integrals	407
A.18	Gradient	408
A.19	Grad Operator	411
A.20	Divergence	412
A.21	Laplacian Operator	415
A.22	Curl	417
A.23	Curvilinear Coordinates	420
A.24	Useful Vector Identities	422

Preface

Physics is an integrated and substantive body of knowledge regarding the nature of the universe that is based on experimental observations. Physics is ultimately expressed as a mathematical model that is capable of both explaining and predicting the behaviors of objects in the natural world. The type of physics taught to undergraduates in universities, and other institutes of tertiary education, has five main components. These components are Newtonian dynamics, classical electromagnetism, special relativity, quantum mechanics, and thermal physics. In universities, undergraduate physics is generally taught at three levels. At the elementary level, students are introduced to the fundamental concepts of Newtonian dynamics, classical electromagnetism, and thermal physics. At the intermediate level, Newtonian dynamics, classical electromagnetism, and thermal physics are fleshed out as relatively coherent theories, and the fundamental concepts of special relativity and quantum mechanics are introduced. At the advanced level, all five components of physics are further developed to their final forms, with the addition of greater abstraction and more advanced mathematics. This course is devoted to intermediate-level physics. The purpose of the course is to present the five components of undergraduate physics as self-consistent and coherent theories. The main emphasis of the presentation is to obtain as many predictions regarding the nature of the physical world as possible while keeping the level of mathematical analysis as low as possible. It turns out that this task is easier to achieve in some areas of physics than in others. In particular, it is not possible to develop a coherent picture of classical electromagnetism without resorting to the sophisticated mathematics of vector field theory.

This course is based on the author's recollection of the first-year survey course, known as Physics Part 1A, that was taught at Cambridge University (U.K.) in the early 1980s. The aim of the course was to bridge the difficult gap between A-level physics and university physics, and also to introduce new concepts in special relativity, quantum mechanics, and thermal physics. For U.S. students, the course bridges the problematic gap between the standard two introductory college physics courses (mechanics/heat and electromagnetism/optics) and upper division physics courses. The course assumes a basic knowledge of physics, trigonometry, algebra, and calculus. The vector algebra and calculus needed to understand the course material is summarized in Appendix A.

Chapter 1

Newtonian Dynamics

1.1 Units

1.1.1 Mks Units

The first principle of any exact physical science is measurement. In Newtonian dynamics, there are three fundamental quantities that are subject to measurement:

1. Intervals in space; that is, *length*.
2. Quantities of inertia, or *inertial mass*, possessed by various bodies.
3. Intervals in *time*.

Any other type of measurement in Newtonian dynamics can (effectively) be reduced to some combination of measurements of these three quantities.

Each of the three fundamental quantities—length, mass, and time—is measured with respect to some convenient standard. The system of units currently used by most scientists and engineers is called the *mks system*—after the first initials of the names of the units of length, mass, and time, respectively, in this system. That is, the *meter*, the *kilogram*, and the *second*.

The mks unit of length is the *meter* (symbol m). The meter was formerly the distance between two scratches on a platinum-iridium alloy bar kept at the International Bureau of Weights and Measures in Sèvres, France, but is now defined as the distance travelled by light in vacuum in $1/299792458$ seconds.

The mks unit of mass is the *kilogram* (symbol kg). The kilogram was formally defined as the mass of a platinum-iridium alloy cylinder kept at the International Bureau of Weights and Measures in Sèvres, France, but is now defined in such a manner as to make Planck's constant take the value $6.626\,070\,15 \times 10^{-34}$ when expressed in mks units.

The mks unit of time is the *second* (symbol s). The second was formerly defined in terms of the Earth's rotation, but is now defined as the time required for 9 192 631 770 complete oscillations associated with the transition between the two hyperfine levels of the ground state of the isotope Cesium 133.

In addition to the three fundamental quantities, Newtonian dynamics also deals with derived quantities, such as velocity, acceleration, momentum, angular momentum, et cetera. Each of these derived quantities can be reduced to some particular combination of length, mass, and time. The mks units of these derived quantities are, therefore, the corresponding combinations of the mks units of length, mass, and time. For instance, a velocity can be reduced to a length divided by a time. Hence, the mks units of velocity are meters per second:

$$[v] = \frac{[L]}{[T]} = \text{m s}^{-1}. \quad (1.1)$$

Here, v stands for a velocity, L for a length, and T for a time, whereas the operator $[\dots]$ represents the units, or dimensions, of the quantity contained within the brackets. Momentum can be reduced to a mass multiplied by a velocity. Hence, the mks units of momentum are kilogram-meters per second:

$$[p] = [M][v] = \frac{[M][L]}{[T]} = \text{kg m s}^{-1}. \quad (1.2)$$

Here, p stands for a momentum, and M for a mass. In this manner, the mks units of all derived quantities appearing in Newtonian dynamics can easily be obtained.

Some combinations of meters, kilograms, and seconds occur so often in physics that they have been given special nicknames. Such combinations include the newton, which is the mks unit of force, and the joule, which is the mks unit of energy. These so-called derived units are listed in Table 1.1.

1.1.2 Standard Prefixes

Mks units are specifically designed to conveniently describe those motions that occur in everyday life. Unfortunately, mks units tend to become rather unwieldy when dealing with motions on very small scales (e.g., the motions of molecules) or on very large scales (e.g., the motions of stars in the Milky Way galaxy). In order to help cope with this problem, a set of standard prefixes has been devised that allow the mks units of length, mass, and time to be modified so as to deal more easily with very small and very large quantities. These prefixes are specified in Table 1.2. Thus, a *kilometer* (km) represents 10^3 m, a *nanometer* (nm) represents 10^{-9} m, and a *femtosecond* (fs) represents 10^{-15} s. The standard prefixes can also be used to modify the units of derived quantities.

Physical Quantity	Derived Unit	Abbreviation	Mks Equivalent
Force	newton	N	m kg s^{-2}
Energy	joule	J	$\text{m}^2 \text{kg s}^{-2}$
Power	watt	W	$\text{m}^2 \text{kg s}^{-3}$
Pressure	pascal	Pa	$\text{m}^{-1} \text{kg s}^{-2}$

Table 1.1: Derived units.

Factor	Prefix	Symbol	Factor	Prefix	Symbol
10^{18}	exa-	E	10^{-1}	deci-	d
10^{15}	peta-	P	10^{-2}	centi-	c
10^{12}	tera-	T	10^{-3}	milli-	m
10^9	giga-	G	10^{-6}	micro-	μ
10^6	mega-	M	10^{-9}	nano-	n
10^3	kilo-	k	10^{-12}	pico-	p
10^2	hecto-	h	10^{-15}	femto-	f
10^1	deka-	da	10^{-18}	atto-	a

Table 1.2: Standard prefixes.

1.1.3 Other Units

The mks system is not the only system of units in existence. Unfortunately, the obsolete cgs (centimeter-gram-second) system, and the even more obsolete fps (foot-pound-second) system, are still in use today, although their continued employment is now strongly discouraged in science and engineering. Conversion between different systems of units is, in principle, perfectly straightforward, but, in practice, a frequent source of error. Witness, for example, the loss of the Mars Climate Orbiter in 1999 (CE) because the Lockheed Martin engineers who designed and built its rocket engine used fps units whereas the NASA mission controllers employed mks units. Table 1.3 specifies the various conversion factors between mks, cgs, and fps units. Note that a pound is a unit of force, rather than mass. Additional non-standard units of length include the inch (1 ft = 12 in), the yard (1 ya = 3 ft), and the mile (1 mi = 5 280 ft). Additional non-standard units of mass include the ton (in the U.S., 1 ton = 2 000 lb; in the U.K., 1 ton = 2 240 lb), and the metric ton (1 tonne = 1 000 kg). Finally, additional non-standard units of time include the minute (1 min = 60 s), the hour (1 hr = 60 min), the (solar) day (1 da = 24 hr), and the (Julian) year (1 yr = 365.25 da).

1 cm	=	10^{-2} m
1 g	=	10^{-3} kg
1 ft	=	0.3048 m
1 lb	=	$4.448 \text{ kg m s}^{-2}$
1 slug	=	14.59 kg

Table 1.3: Conversion factors between the mks, cgs, and fps systems of units. Here, g, ft, and lb are the abbreviations for gram, foot, and pound, respectively.

1.1.4 Dimensional Analysis

As we have already seen, length, mass, and time are three fundamentally different entities that are measured in terms of three completely independent units. It, therefore, makes no sense for a prospective law of physics to express an equality between (say) a length and a mass. In other words, the prospective physical law,

$$m = l, \quad (1.3)$$

where m is a mass and l is a length, cannot possibly be correct. One easy way of seeing that Equation (1.3) is invalid (as a law of physics) is to note that this equation is dependent on the adopted system of units. That is, if $m = l$ in mks units then $m \neq l$ in fps units, because the conversion factors which must be applied to the left- and right-hand sides of the equation differ. Physicists hold very strongly to the maxim that the laws of physics possess objective reality. In other words, the laws of physics are equivalent for all observers. One immediate consequence of this maxim is that a law of physics must take the same form in all possible systems of units that a prospective observer might choose to employ (because the choice of units is arbitrary, and has nothing to do with physical reality). The only way in which this can be the case is if all laws of physics are dimensionally consistent. In other words, the quantities on the left- and right-hand sides of the equality sign in any given law of physics must have the same dimensions (i.e., the same combinations of length, mass, and time). A dimensionally consistent equation naturally takes the same form in all possible systems of units, because the same conversion factors are applied to both sides of the equation when transforming from one system to another.

As an example, let us consider what is probably the most famous equation in physics; that is, Einstein's mass-energy relation,

$$E = m c^2. \quad (1.4)$$

(See Section 3.3.4.) Here, E is the energy of a body, m is its mass, and c is the speed of light in vacuum. The dimensions of energy are $[M][L]^2/[T]^2$, and the dimensions of speed are $[L]/[T]$. Hence, the dimensions of the left-hand side are $[M][L]^2/[T]^2$, whereas the dimensions of the right-hand side are $[M]([L]/[T])^2 = [M][L]^2/[T]^2$. It follows that Equation (1.4) is indeed dimensionally consistent. Thus, $E = m c^2$ holds good in mks units, in cgs units, in fps units, and in any other sensible set of units. Had Einstein proposed $E = m c$, or $E = m c^3$ then his error would have been immediately apparent to other physicists, because these prospective laws are not dimensionally consistent. In fact, $E = m c^2$ represents the only simple, dimensionally consistent way of combining an energy, a mass, and the velocity of light in a law of physics.

The last comment leads naturally to the subject of *dimensional analysis*. That is, the use of the idea of dimensional consistency to guess the forms of simple laws of physics.

1.1.5 Deriving Physical Relationships

Consider a viscous fluid flowing through a circular pipe. The volume rate of fluid flow through the pipe, Q , might plausibly depend on the radius of the pipe, a , the viscosity of the fluid, η , and the pressure gradient along the pipe, $\Delta p/l$. Here, Δp is the pressure difference between the two ends

of the pipe, and l is the length of the pipe. Let us guess that

$$Q = A a^x \eta^y \left(\frac{\Delta p}{l} \right)^z, \quad (1.5)$$

where x , y , and z are, as yet, unknown exponents, and A is a dimensionless constant. Now, the dimensions of Q are $[L]^3/[T]$, the dimensions of a are $[L]$, the dimensions of η are $[M]/([L][T])$, and the dimensions of $\Delta p/l$ are $\{([M][L]/[T]^2)/[L]^2\}/[L] = [M]/([L]^2[T]^2)$. Thus, equating the dimensions of the left- and right-hand sides of the previous equation, we obtain

$$\frac{[L]^3}{[T]} = [L]^x \left(\frac{[M]}{[L][T]} \right)^y \left(\frac{[M]}{[L]^2[T]^2} \right)^z. \quad (1.6)$$

Now, if Equation (1.5) is to be dimensionally consistent then we can separately equate the exponents of length, mass, and time in the previous expression. Equating the exponents of $[L]$, we obtain

$$3 = x - y - 2z. \quad (1.7)$$

Equating the exponents of $[M]$, we get

$$0 = y + z. \quad (1.8)$$

Finally, equating the exponents of $[T]$, we obtain

$$-1 = -y - 2z. \quad (1.9)$$

It is easily seen that $x = 4$, $y = -1$, and $z = 1$. Hence, we deduce that

$$Q = A \frac{a^4}{\eta} \left(\frac{\Delta p}{l} \right). \quad (1.10)$$

1.1.6 Scaling Laws

Suppose that a special effects studio wants to film a scene in which the Leaning Tower of Pisa topples to the ground. In order to achieve this goal, the studio might make a scale model of the tower, which is (say) 1 m tall, and then film the model falling over. The only problem is that the resulting footage would look completely unrealistic because the model tower would fall over too quickly. The studio could easily fix this problem by slowing the film down. But, by what factor should the film be slowed down in order to make it look realistic?

Although, at this stage, we do not know how to apply the laws of physics to the problem of a tower falling over, we can, at least, make some educated guesses as to the factors upon which the time, t_f , required for this process to occur depends. In fact, it seems reasonable to suppose that t_f depends principally on the mass of the tower, m , the height of the tower, h , and the acceleration due to gravity, g . In other words,

$$t_f = C m^x h^y g^z, \quad (1.11)$$

where C is a dimensionless constant, and x , y , and z are unknown exponents. The exponents x , y , and z can be determined by the requirement that the previous equation be dimensionally consistent.

Incidentally, the dimensions of an acceleration are $[L]/[T]^2$. Hence, equating the dimensions of both sides of Equation (1.11), we obtain

$$[T] = [M]^x [L]^y \left(\frac{[L]}{[T]^2} \right)^z. \quad (1.12)$$

We can now compare the exponents of $[L]$, $[M]$, and $[T]$ on either side of the previous expression. These exponents must all match in order for Equation (1.11) to be dimensionally consistent. Thus,

$$0 = y + z, \quad (1.13)$$

$$0 = x, \quad (1.14)$$

$$1 = -2z. \quad (1.15)$$

It immediately follows that $x = 0$, $y = 1/2$, and $z = -1/2$. Hence,

$$t_f = C \sqrt{\frac{h}{g}}. \quad (1.16)$$

Now, the actual tower of Pisa is approximately 100 m tall. It follows that because $t_f \propto \sqrt{h}$ (g is the same for both the real and the model tower) the 1 m high model tower would fall over a factor of $\sqrt{100/1} = 10$ times faster than the real tower. Thus, the film must be slowed down by a factor of 10 in order to make it look realistic.

1.2 Newton's Laws of Motion

1.2.1 Introduction

Newton's laws of motion were first enunciated by Sir Isaac Newton in a work entitled *Philosophiae Naturalis Principia Mathematica* (Mathematical Principles of Natural Philosophy). This work, which was first published in 1687 (CE), is nowadays more commonly referred to as the *Principia*.

1.2.2 Newton's First Law of Motion

Newton's first law of motion is

Every body continues in its state of rest, or uniform motion in a straight line, unless compelled to change that state by forces impressed upon it.

Newton's first law of motion states that a body subject to zero net force does not accelerate (i.e., it moves in a straight line at a constant speed). However, this law is only valid in special frames of reference known as *inertial frames*. In fact, we can think of Newton's first law as the definition of an inertial frame. Namely, an inertial reference frame is one in which a body subject to zero net force does not accelerate. There are an infinite number of different inertial reference frames

all moving at constant velocities with respect to one another. (See Section 1.5.4.) It is impossible to identify an exact inertial reference frame. The best approximation to such a frame is the so-called International Celestial Reference System (ICRF), whose origin is the center of mass of the solar system, and whose coordinate axes are defined with respect to extremely distant point radio sources (mostly quasars) whose positions can be measured to great accuracy via very long baseline interferometry (VLBI).¹

1.2.3 Newton's Second Law of Motion

Newton's second law of motion is

The change of motion (i.e., momentum) of a body is proportional to the force impressed upon it, and is made in the direction of the straight line in which the force is impressed.

As before, Newton's second law is only valid in an inertial reference frame. Suppose that the body in question has a *mass* m , a *displacement* (from an arbitrary stationary point that forms the origin of a Cartesian coordinate system that we have set up in our inertial reference frame) \mathbf{r} , an instantaneous *velocity* $\mathbf{v} = d\mathbf{r}/dt$, and is subject to a *force* \mathbf{f} . Here, t denotes time. Newton's second law of motion states that

$$\mathbf{f} = \frac{d\mathbf{p}}{dt} = \frac{d(m\mathbf{v})}{dt}, \quad (1.17)$$

where

$$\mathbf{p} = m\mathbf{v} \quad (1.18)$$

is the body's *linear momentum*.

If the mass of the body is assumed to be constant then Equation (1.17) reduces to

$$\mathbf{f} = m \frac{d\mathbf{v}}{dt} = m\mathbf{a}, \quad (1.19)$$

where $\mathbf{a} = d\mathbf{v}/dt$ is the body's instantaneous *acceleration*. Note that the mass that appears in the previous equation is a measure of the reluctance of the body to deviate from its preferred state of uniform motion in a straight line due to the action of a force. This type of mass is known as *inertial mass*. However, another type of mass occurs in nature. A body situated in a gravitational field whose local acceleration is \mathbf{g} is subject to a gravitational force

$$\mathbf{f} = m\mathbf{g}. \quad (1.20)$$

(See Section 1.8.1.) The mass that appears in the previous equation is a sort of gravitational charge (i.e., it is analogous to the electric charge of a particle in an electric field). This type of mass is known as *gravitational mass*. It is an observational fact that inertial mass is proportional to gravitational mass for all bodies in the universe. In fact, the conventional system of units used in physics is set up in such a manner that inertial mass is equal to gravitational mass. Nevertheless, it

¹P. Charlot, et al., *Astronomy and Astrophysics* **644**, A159 (2020).

is important to understand that these two types of masses measure different physical properties of a given body.

Incidentally, the reason that inertial mass is proportional to gravitational mass was not explained until 1916, when Albert Einstein proposed his *general theory of relativity*. According to this theory, inertial mass is proportional to gravitational mass because it is impossible to distinguish experimentally between a gravitational acceleration and a fictitious acceleration due to motion observed in a non-inertial reference frame. (See Section 1.5.4.)

Acceleration is a *vector* (i.e., it transforms under rotation of the coordinate axes in an analogous manner to a displacement), whereas mass is a *scalar* (i.e., it is invariant under rotation of the coordinate axes). (See Section A.5.) Thus, it follows from Equation (1.19) that force must be a vector. (Otherwise, the form of Newton's second law would depend unphysically on the arbitrary orientation of the coordinate axes.) One consequence of force being a vector is that two forces, \mathbf{f}_1 and \mathbf{f}_2 , both acting on a given body, have the same effect as a single force, $\mathbf{f} = \mathbf{f}_1 + \mathbf{f}_2$, acting on the same body, where the summation is performed according to the laws of vector addition. (See Section A.3.) Likewise, a single force, \mathbf{f} , acting at on a given body has the same effect as two forces, \mathbf{f}_1 and \mathbf{f}_2 , acting on the same body, provided that $\mathbf{f}_1 + \mathbf{f}_2 = \mathbf{f}$. This method of combining and splitting forces is known as the *resolution of forces*, and lies at the heart of many calculations in classical dynamics.

1.2.4 Newton's Third Law of Motion

Newton's third law of motion is

To every action there is always opposed an equal reaction; or, the mutual actions of two bodies upon each other are always equal and directed to contrary parts.

Consider a dynamical system consisting of two bodies, labelled 1 and 2. Let body 1 exert a force \mathbf{f}_{21} on body 2, and let body 2 exert a force \mathbf{f}_{12} on body 1. According to Newton's third law of motion,

$$\mathbf{f}_{12} = -\mathbf{f}_{21}. \quad (1.21)$$

In other words, the two forces are equal and opposite. In Newtonian language, one of the forces is the *action*, and the other is the *reaction*. Thus, action and reaction are always equal and opposite. Newton's third law holds irrespective of the nature of the forces acting between the two bodies. One corollary of this law is that a body cannot exert a force on itself. Another corollary is that all (non-fictitious) forces in the universe have corresponding reactions.

It should be noted that Newton's third law implies *action at a distance*. In other words, if the force that body 1 exerts on body 2 suddenly changes then Newton's third law demands that there must be an immediate change in the force that body 2 exerts on body 1. Moreover, this must be the case irrespective of the distance between the two bodies. However, we know that Einstein's special theory of relativity forbids information from traveling through the universe faster than the speed of light in vacuum. (See Section 3.2.10.) Hence, action at a distance is also forbidden. In other words, if the force that body 1 exerts on body 2 suddenly changes then there must be a time delay, which is at least as long as it takes a light ray to propagate between the two bodies, before

the force that body 2 exerts on body 1 can respond. Of course, this means that Newton's third law is not, strictly speaking, correct. However, as long as we restrict our investigations to the motions of dynamical systems on timescales that are long compared to the time required for light-rays to traverse these systems, Newton's third law can be regarded as being approximately correct.

1.3 Motion of Single Particle

1.3.1 Impulse

Consider the motion of a single particle (i.e., a body of negligible spatial extent). Newton's second law of motion, (1.19), can be written

$$\mathbf{f} = m \frac{d\mathbf{v}}{dt}, \quad (1.22)$$

which implies that

$$\mathbf{f} dt = m d\mathbf{v}. \quad (1.23)$$

Suppose that the particle in question has an instantaneous velocity \mathbf{v}_1 at an initial time t_1 , and an instantaneous velocity \mathbf{v}_2 at a final time t_2 . Integrating the previous equation between the initial and the final time, we obtain

$$\int_{t_1}^{t_2} \mathbf{f}(t) dt = m \int_{\mathbf{v}_1}^{\mathbf{v}_2} d\mathbf{v} = m(\mathbf{v}_2 - \mathbf{v}_1), \quad (1.24)$$

where we have taken into account the fact that the force \mathbf{f} is, in general, a function of time. The quantity

$$\mathbf{I} = \int_{t_1}^{t_2} \mathbf{f}(t) dt \quad (1.25)$$

is known as *impulse*, and is essentially the 'area' under the $\mathbf{f}(t)$ curve between times t_1 and t_2 . It is clear from the previous two equations that

$$\mathbf{I} = m \Delta\mathbf{v}, \quad (1.26)$$

where $\Delta\mathbf{v} = \mathbf{v}_2 - \mathbf{v}_1$ is the change in the particle's velocity between the initial and the final times. Equation (1.18) can be combined with the previous two equations to give

$$\mathbf{I} = \Delta\mathbf{p}. \quad (1.27)$$

In other words, the net impulse acting on a particle between an initial and a final time is equal to the change in the momentum of the particle between the same two times.

1.3.2 Work

Suppose that a particle subject to a force \mathbf{f} undergoes an infinitesimal displacement $d\mathbf{r}$. The net work that the force does on the particle (i.e., the net energy transferred to the body by the force) is

$$dW = \mathbf{f} \cdot d\mathbf{r}. \quad (1.28)$$

(See Section A.6.) In other words, the work is the product of the displacement and the component of the force parallel to the displacement. It follows from Equation (1.22) that

$$dW = m \frac{d\mathbf{v}}{dt} \cdot d\mathbf{r}. \quad (1.29)$$

However, $d\mathbf{r} = \mathbf{v} dt$, so we obtain

$$dW = m \frac{d\mathbf{v}}{dt} \cdot \mathbf{v} dt = m \mathbf{v} \cdot d\mathbf{v}. \quad (1.30)$$

Furthermore,

$$\mathbf{v} \cdot d\mathbf{v} = v_x dv_x + v_y dv_y + v_z dv_z = \frac{1}{2} d(v_x^2) + \frac{1}{2} d(v_y^2) + \frac{1}{2} d(v_z^2) = \frac{1}{2} d(v^2), \quad (1.31)$$

where $v = |\mathbf{v}| = (v_x^2 + v_y^2 + v_z^2)^{1/2}$ is the particle's speed. It follows from the previous two equations that

$$dW = dK, \quad (1.32)$$

where

$$K = \frac{1}{2} m v^2. \quad (1.33)$$

Here, K is known as *kinetic energy*, and is the energy that the particle possesses by virtue of its motion. Equation (1.32) can be integrated to give the *work-energy theorem*,

$$W = \Delta K. \quad (1.34)$$

According to this theorem, the net work done by the force acting on the particle in a given time interval is equal to the change in the particle's kinetic energy during the same time interval.

Suppose that the force is a function of the particle's displacement, \mathbf{r} . If the particle moves from point A to point B along any path then Equations (1.28) and (1.34) imply that

$$W = \int_{\mathbf{r}_A}^{\mathbf{r}_B} \mathbf{f} \cdot d\mathbf{r} = K_B - K_A, \quad (1.35)$$

where \mathbf{r}_A denotes the displacement of point A , et cetera, K_A is the kinetic energy at point A , et cetera, and $d\mathbf{r}$ is an element of the path. (See Section A.14.)

1.3.3 Conservative Forces

Suppose, again, that a particle is acted upon by a force $\mathbf{f}(\mathbf{r})$ that is a function of the particle's displacement, \mathbf{r} . Suppose that the body travels from point A to point B along some particular path, labelled 1. The net work done on the particle is

$$W_1 = \left(\int_{\mathbf{r}_A}^{\mathbf{r}_B} \mathbf{f} \cdot d\mathbf{r} \right)_{\text{path 1}}, \quad (1.36)$$

where $d\mathbf{r}$ is an element of the path. (See Section A.14.) Suppose, now, that the particle travels between the same two points along a different path, labelled 2. The net work done on the particle is

$$W_2 = \left(\int_{\mathbf{r}_A}^{\mathbf{r}_B} \mathbf{f} \cdot d\mathbf{r} \right)_{\text{path 2}}. \quad (1.37)$$

There are two types of forces in the universe. *Conservative* forces are such that

$$W_1 = W_2 \quad (1.38)$$

irrespective of the locations of points A and B , and the nature of paths 1 and 2. (See Section A.18.) In other words, a conservative force is such that the net work done on a particle moving between two points is independent of the path taken between the two points. Gravity is an example of a conservative force. On the other hand, *non-conservative* forces are such that net work done on a particle moving between two points depends on the path taken between the two points. Friction is an example of a non-conservative force.

Suppose that the particle is acted on by a conservative force and moves from point A to point B along path 1, and then from point B to point A along path 2. In other words, the particle moves in a closed loop. The net work done on the particle is

$$\begin{aligned} W &= \left(\int_{\mathbf{r}_A}^{\mathbf{r}_B} \mathbf{f} \cdot d\mathbf{r} \right)_{\text{path 1}} + \left(\int_{\mathbf{r}_B}^{\mathbf{r}_A} \mathbf{f} \cdot d\mathbf{r} \right)_{\text{path 2}} \\ &= \left(\int_{\mathbf{r}_A}^{\mathbf{r}_B} \mathbf{f} \cdot d\mathbf{r} \right)_{\text{path 1}} - \left(\int_{\mathbf{r}_A}^{\mathbf{r}_B} \mathbf{f} \cdot d\mathbf{r} \right)_{\text{path 2}} \\ &= W_1 - W_2 = 0, \end{aligned} \quad (1.39)$$

where use has been made of the previous three equations. Thus, we conclude that

$$\oint \mathbf{f} \cdot d\mathbf{r} = 0. \quad (1.40)$$

(See Section A.18.) In other words, if a particle subject to a conservative force moves in a closed loop then zero net work is done on the particle.

1.3.4 Potential Energy

Consider a particle subject to a conservative force. Let O be the origin of our coordinate system (i.e., the point whose displacement is $\mathbf{0}$), and let P be a general point whose displacement is \mathbf{r} . We can define the function

$$U(\mathbf{r}) = - \int_{\mathbf{0}}^{\mathbf{r}} \mathbf{f}(\mathbf{r}') \cdot d\mathbf{r}'. \quad (1.41)$$

The fact that the force is conservative ensures that this function has a unique value at each point in space. On the other hand, if the force were non-conservative then the function would be ill-defined, because there are an infinite number of different paths linking points O and P , and each path would yield a different value of the integral on the right-hand side of the previous equation. The quantity U is known as *potential energy*, and is the energy that the particle possesses by virtue of its position. Obviously, it only makes sense to associate potential energy with a conservative force. Note that the fact that the position of the origin of our coordinate system is arbitrary implies that potential energy is undefined to an arbitrary additive constant. In other words, only differences in potential energies are physically meaningful.

Suppose that the particle moves from point \mathbf{r} to point $\mathbf{r} + d\mathbf{r}$. The associated change in the particle's potential energy is

$$dU = -\mathbf{f} \cdot d\mathbf{r} = -f_x dx - f_y dy - f_z dz. \quad (1.42)$$

Suppose that $dy = dz = 0$. We can write

$$f_x = - \left(\frac{dU}{dx} \right)_{\text{constant } y, z} = - \frac{\partial U}{\partial x}. \quad (1.43)$$

Similar arguments yield

$$f_y = - \frac{\partial U}{\partial y}, \quad (1.44)$$

$$f_z = - \frac{\partial U}{\partial z}. \quad (1.45)$$

Hence, we deduce that

$$\mathbf{f} = - \frac{\partial U}{\partial x} \mathbf{e}_x - \frac{\partial U}{\partial y} \mathbf{e}_y - \frac{\partial U}{\partial z} \mathbf{e}_z, \quad (1.46)$$

where \mathbf{e}_x is a unit vector parallel to the x -axis, et cetera. (See Section A.4.) The previous equation can be written more succinctly as

$$\mathbf{f} = -\nabla U. \quad (1.47)$$

(See Section A.19.) In other words, a particle moving in a conservative force field experiences a force that is equal to minus the *gradient* of the potential energy.

1.3.5 Energy Conservation

Consider a particle moving in a conservative force field. Suppose that the particle moves from point A to point B along some particular path. According to Equation (1.35),

$$\int_{\mathbf{r}_A}^{\mathbf{r}_B} \mathbf{f} \cdot d\mathbf{r} = K_B - K_A. \quad (1.48)$$

However, Equation (1.41) implies that

$$\int_{\mathbf{r}_A}^{\mathbf{r}_B} \mathbf{f} \cdot d\mathbf{r} = - \int_0^{\mathbf{r}_A} \mathbf{f} \cdot d\mathbf{r} + \int_0^{\mathbf{r}_B} \mathbf{f} \cdot d\mathbf{r} = U_A - U_B, \quad (1.49)$$

where U_A is the potential energy at point A , et cetera. The previous two equations yield

$$K_A + U_A = K_B + U_B. \quad (1.50)$$

Thus, if we define the total energy, E , of the particle as the sum of its kinetic and potential energies,

$$E = K + U, \quad (1.51)$$

then we deduce that E is a constant of the motion. In other words, the total energy of a particle moving in a conservative force field is a conserved quantity.

1.3.6 Energy Diagrams

Consider a moving in the x -direction, say, under the action of some x -directed force, $f(x)$. Suppose that $f(x)$ is a conservative force; for instance, gravity. In this case, according to Equation (1.47), we can write

$$f(x) = -\frac{dU(x)}{dx}, \quad (1.52)$$

where $U(x)$ is the potential energy of the particle at position x .

Let the curve $U(x)$ take the form shown in Figure 1.1. For instance, this curve might represent the gravitational potential energy of a cyclist freewheeling in a hilly region. Observe that we have set the potential energy at infinity to zero (which we are generally free to do, because potential energy is undefined to an arbitrary additive constant). This is a fairly common convention. What can we deduce about the motion of the particle in this potential?

We know that the total energy, E —which is the sum of the kinetic energy, K , and the potential energy, U —is a constant of the motion. [See Equation (1.51).] Hence, we can write

$$K(x) = E - U(x). \quad (1.53)$$

However, we also know that a kinetic energy can never be negative [because $K = (1/2)mv^2$, and neither m nor v^2 can be negative]. Hence, the previous expression tells us that the particle's motion is restricted to the region (or regions) in which the potential energy curve $U(x)$ falls below the value E . This idea is illustrated in Figure 1.1. Suppose that the total energy of the system is E_0 .

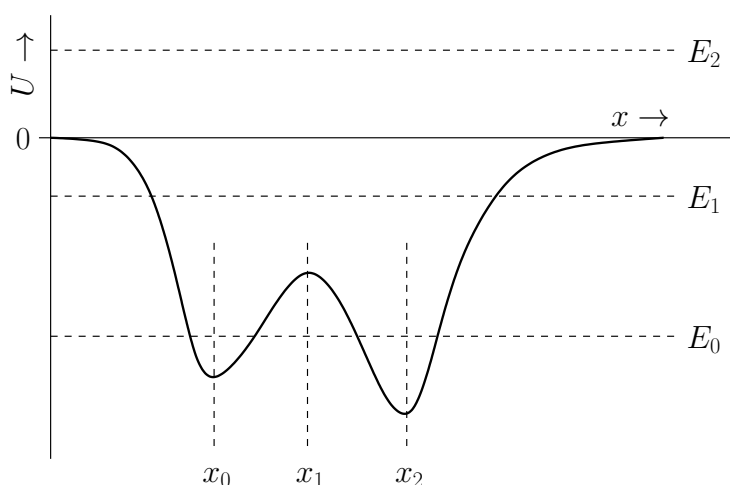


Figure 1.1: A potential energy curve.

It is clear, from the figure, that the particle is trapped inside one or other of the two dips in the potential; these dips are generally referred to as *potential wells*. Suppose that we now raise the energy to E_1 . In this case, the particle is free to enter or leave each of the potential wells, but its motion is still bounded to some extent, because it clearly cannot move off to infinity. Finally, let us raise the energy to E_2 . Now the particle is unbounded; that is, it can move off to infinity. In conservative systems in which it makes sense to adopt the convention that the potential energy at infinity is zero, bounded systems are characterized by $E < 0$, whereas unbounded systems are characterized by $E > 0$.

The previous discussion suggests that the motion of a particle moving in a potential generally becomes less bounded as the total energy E of the system increases. Conversely, we would expect the motion to become more bounded as E decreases. In fact, if the energy becomes sufficiently small then it appears likely that the system will settle down in some equilibrium state in which the particle is stationary. Let us try to identify any prospective equilibrium states in Figure 1.1. If the particle remains stationary then it must be subject to zero force (otherwise it would accelerate). Hence, according to Equation (1.52), an equilibrium state is characterized by

$$\frac{dU}{dx} = 0. \quad (1.54)$$

In other words, an equilibrium state corresponds to either a maximum or a minimum of the potential energy curve $U(x)$. It can be seen that the $U(x)$ curve shown in Figure 1.1 has three associated equilibrium states located at $x = x_0$, $x = x_1$, and $x = x_2$.

Let us now make a distinction between stable equilibrium points and unstable equilibrium points. When the particle is slightly displaced from a *stable* equilibrium point then the resultant force f acting on it must always be such as to return it to this point. In other words, if $x = x_0$ is an equilibrium point then we require

$$\left. \frac{df}{dx} \right|_{x=x_0} < 0 \quad (1.55)$$

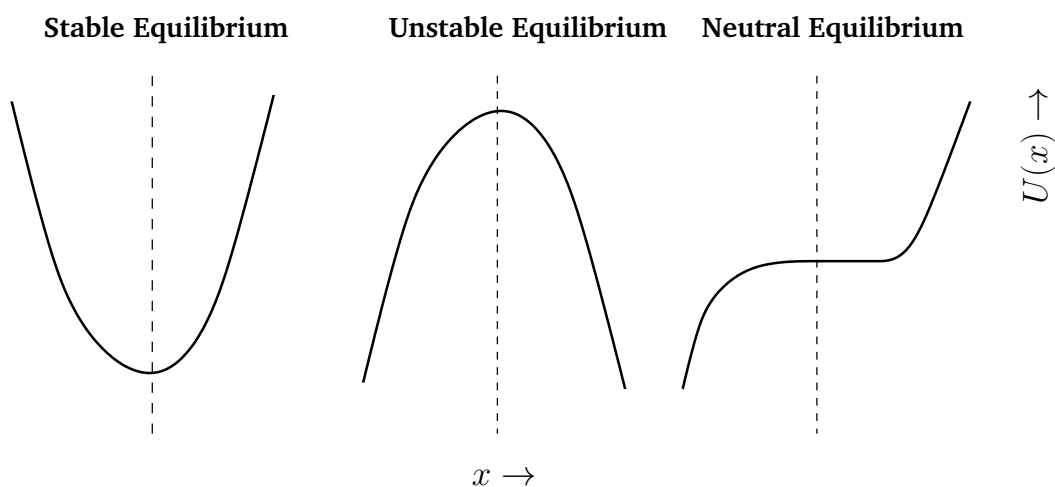


Figure 1.2: Different types of equilibrium point.

for stability; that is, if the particle is displaced to the right, so that $x - x_0 > 0$, then the force must act to the left, so that $f < 0$, and vice versa. Likewise, if

$$\left. \frac{df}{dx} \right|_{x=x_0} > 0 \quad (1.56)$$

then the equilibrium point $x = x_0$ is *unstable*. It follows, from Equation (1.52), that stable equilibrium points are characterized by

$$\frac{d^2U}{dx^2} > 0. \quad (1.57)$$

In other words, a stable equilibrium point corresponds to a minimum of the potential energy curve $U(x)$. Likewise, an unstable equilibrium point corresponds to a maximum of the $U(x)$ curve. Hence, we conclude that, in Figure 1.1, $x = x_0$ and $x = x_2$ are stable equilibrium points, whereas $x = x_1$ is an unstable equilibrium point. Of course, this makes perfect sense if we think of $U(x)$ as a gravitational potential energy curve, so that U is directly proportional to height. In this case, all we are saying is that it is easy to confine a low energy mass at the bottom of a valley, but very difficult to balance the same mass on the top of a hill (because any slight displacement of the mass will cause it to slide down the hill). Note, finally, that if

$$\frac{dU}{dx} = \frac{d^2U}{dx^2} = 0 \quad (1.58)$$

at any point (or in any region) then we have what is known as a *neutral* equilibrium point. We can move the particle slightly away from such a point and it will still remain in equilibrium (i.e., it will neither attempt to return to its initial state, nor will it continue to move). A neutral equilibrium point corresponds to a flat spot in a $U(x)$ curve. See Figure 1.2.

The equation of motion of a particle moving in one dimension under the action of a conservative force is, in principle, integrable. Because $K = (1/2)mv^2$, the energy conservation equation (1.53)

can be rearranged to give

$$v = \pm \left(\frac{2[E - U(x)]}{m} \right)^{1/2}, \quad (1.59)$$

where the \pm signs correspond to motion to the left and to the right, respectively. However, because $v = dx/dt$, this expression can be integrated to give

$$t = \pm \left(\frac{m}{2E} \right)^{1/2} \int_{x_0}^x \frac{dx'}{\sqrt{1 - U(x')/E}}, \quad (1.60)$$

where $x(t = 0) = x_0$. For sufficiently simple potential functions, $U(x)$, the previous equation can be solved to give x as a function of t . For instance, if $U = (1/2)kx^2$, $x_0 = 0$, and the plus sign is chosen, then

$$t = \left(\frac{m}{k} \right)^{1/2} \int_0^{(k/2E)^{1/2}x} \frac{dy}{\sqrt{1 - y^2}} = \left(\frac{m}{k} \right)^{1/2} \sin^{-1} \left(\left[\frac{k}{2E} \right]^{1/2} x \right), \quad (1.61)$$

which can be inverted to give

$$x = a \sin(\omega t), \quad (1.62)$$

where $a = \sqrt{2E/k}$ and $\omega = \sqrt{k/m}$. This type of motion is known as *simple harmonic motion*. Note that the particle reverses direction each time it reaches one of the so-called *turning points* ($x = \pm a$) at which $U = E$ (and, so $K = 0$). This analysis suggests that a particle trapped in a general potential well exhibits oscillatory motion between the turning points.

1.4 Motion of System of Many Particles

1.4.1 Equations of Motion

Consider a dynamical system consisting of N particles. Let particle i have mass m_i , displacement \mathbf{r}_i , and velocity $\mathbf{v}_i = d\mathbf{r}_i/dt$. Suppose that particle i is subject to a force \mathbf{f}_{ij} exerted by particle j . Suppose, in addition, that particle i is subject to an external force (i.e., a force that originates outside the dynamical system) \mathbf{F}_i . Applying Newton's second law of motion to the particle [see Equation (1.19)], we obtain

$$m_i \frac{d\mathbf{v}_i}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{f}_{ij} + \mathbf{F}_i, \quad (1.63)$$

assuming that all of the forces acting on particle i are superposable. (This is reasonable because gravitational and electromagnetic forces are superposable.) Newton's third law of motion, (1.21), can be generalized to give

$$\mathbf{f}_{ij} = -\mathbf{f}_{ji}, \quad (1.64)$$

for all i and j . Note, in particular, that $\mathbf{f}_{ii} = -\mathbf{f}_{ii} = \mathbf{0}$. In other words, particle i cannot exert a force on itself. This accounts for the exclusion of particle i in the sum on the right-hand side of Equation (1.63).

There are N equations of motion of analogous form to Equation (1.63); one for each particle that makes up the system. We can sum all of these equations to give

$$\sum_{i=1,N} m_i \frac{d\mathbf{v}_i}{dt} = \sum_{i=1,N} \sum_{\substack{j=1,N \\ j \neq i}} \mathbf{f}_{ij} + \sum_{i=1,N} \mathbf{F}_i. \quad (1.65)$$

Now, every term, \mathbf{f}_{ij} , appearing in the double sum on the right-hand side of the previous equation, can be paired with another term— \mathbf{f}_{ji} , in this case—that is equal and opposite according to Newton's third law of motion, (1.64). In other words, the terms in the sum all cancel out in pairs. It follows that the previous equation reduces to

$$\sum_{i=1,N} m_i \frac{d\mathbf{v}_i}{dt} = \mathbf{F}, \quad (1.66)$$

where

$$\mathbf{F} = \sum_{i=1,N} \mathbf{F}_i. \quad (1.67)$$

is the net external force acting on the system.

1.4.2 Center of Mass

The center of mass of a dynamical system is an imaginary point whose coordinates are the mass-weighted average of the coordinates of the system's constituent particles. It follows that the displacement of the center of mass is

$$\mathbf{R} = \frac{\sum_{i=1,N} m_i \mathbf{r}_i}{\sum_{i=1,N} m_i}. \quad (1.68)$$

The velocity of the center of mass, which is obtained by differentiating the previous expression with respect to time, is

$$\frac{d\mathbf{R}}{dt} = \frac{1}{M} \sum_{i=1,N} m_i \mathbf{v}_i, \quad (1.69)$$

where

$$M = \sum_{i=1,N} m_i \quad (1.70)$$

is the total mass of the system. Likewise, the acceleration of the center of mass is

$$\frac{d^2\mathbf{R}}{dt^2} = \frac{1}{M} \sum_{i=1,N} m_i \frac{d\mathbf{v}_i}{dt}. \quad (1.71)$$

A comparison of Equations (1.66) and (1.71) reveals that

$$M \frac{d^2\mathbf{R}}{dt^2} = \mathbf{F}. \quad (1.72)$$

We conclude that the center of mass moves like a particle of mass M subject to the net external force, \mathbf{F} , acting on the system. In particular, the motion of the center of mass is completely unaffected by the internal forces that the system's constituent particles exert on one another.

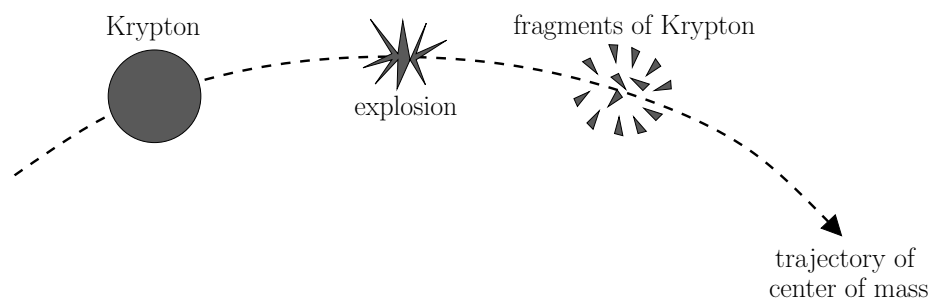


Figure 1.3: The unfortunate history of the planet Krypton.

1.4.3 Explosion of Krypton

As an illustration of the points raised in the previous discussion, let us consider the unfortunate history of the planet Krypton. As is well-known, Krypton—Superman’s home planet—eventually exploded. Note, however, that before, during, and after this explosion, the net external force acting on Krypton, or the fragments of Krypton—namely, the gravitational attraction due to Krypton’s sun—remained the same. In other words, the forces responsible for the explosion can be thought of as large, transitory, internal forces. We conclude that the motion of the center of mass of Krypton, or the fragments of Krypton, was unaffected by the explosion. This follows, from Equation (1.72), because the motion of the center of mass is independent of internal forces. Before the explosion, the planet Krypton presumably executed a standard Keplerian orbit around Krypton’s sun. We conclude that, after the explosion, the fragments of Krypton (or, to be more exact, the center of mass of these fragments) continued to execute exactly the same orbit. See Figure 1.3.

1.4.4 Conservation of Linear Momentum

Suppose that our dynamical system is *isolated*. In other words, the system is not subject to a net external force, so that $\mathbf{F} = \mathbf{0}$. In this case, Equation (1.66) reduces to

$$\sum_{i=1,N} m_i \frac{d\mathbf{v}_i}{dt} = \mathbf{0}. \quad (1.73)$$

However, the linear momentum of the i th particle is

$$\mathbf{p}_i = m_i \mathbf{v}_i. \quad (1.74)$$

Thus, Equation (1.73) yields

$$\sum_{i=1,N} \frac{d\mathbf{p}_i}{dt} = \mathbf{0}, \quad (1.75)$$

or

$$\frac{d\mathbf{P}}{dt} = \mathbf{0}, \quad (1.76)$$

where

$$\mathbf{P} = \sum_{i=1,N} \mathbf{p}_i \quad (1.77)$$

is the total linear momentum of the system. Equation (1.76) implies that the total linear momentum of an isolated dynamical system is a conserved quantity. In other words, the total momentum does not evolve in time.

1.4.5 Conservation of Angular Momentum

Taking the vector product of Equation (1.63) with the displacement \mathbf{r}_i , we obtain

$$m_i \mathbf{r}_i \times \frac{d\mathbf{v}_i}{dt} = \sum_{j=1,N}^{j \neq i} \mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_i \times \mathbf{F}_i. \quad (1.78)$$

Now, the *angular momentum* of particle i about the origin of the coordinate system is defined

$$\mathbf{l}_i = \mathbf{r}_i \times \mathbf{p}_i = m_i \mathbf{r}_i \times \mathbf{v}_i. \quad (1.79)$$

(See Section A.8.) It follows that

$$\frac{d\mathbf{l}_i}{dt} = m_i \mathbf{v}_i \times \mathbf{v}_i + m_i \mathbf{r}_i \times \frac{d\mathbf{v}_i}{dt} = m_i \mathbf{r}_i \times \frac{d\mathbf{v}_i}{dt}. \quad (1.80)$$

(See Sections A.8 and A.12.) Hence, Equation (1.78) yields the following angular equation of motion for the i th particle:

$$\frac{d\mathbf{l}_i}{dt} = \sum_{j=1,N}^{j \neq i} \mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_i \times \mathbf{F}_i. \quad (1.81)$$

There are N angular equations of motion of analogous form to the previous equation; one for each particle that makes up the system. We can sum all of these equations to give

$$\frac{d\mathbf{L}}{dt} = \sum_{i=1,N} \sum_{j=1,N}^{j \neq i} \mathbf{r}_i \times \mathbf{f}_{ij} + \boldsymbol{\tau}, \quad (1.82)$$

where

$$\mathbf{L} = \sum_{i=1,N} \mathbf{l}_i \quad (1.83)$$

is the total angular momentum of the system about the origin of the coordinate system, and

$$\boldsymbol{\tau} = \sum_{i=1,N} \mathbf{r}_i \times \mathbf{F}_i \quad (1.84)$$

is the net external torque acting on the system about the origin of the coordinate system. (See Section A.8.)

Consider the double sum on the right-hand side of Equation (1.82). A general term, $\mathbf{r}_i \times \mathbf{f}_{ij}$, in this sum can always be paired with a matching term, $\mathbf{r}_j \times \mathbf{f}_{ji}$, in which the indices have been swapped. Making use of Equation (1.64), the sum of a general matched pair can be written

$$\mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_j \times \mathbf{f}_{ji} = (\mathbf{r}_i - \mathbf{r}_j) \times \mathbf{f}_{ij}. \quad (1.85)$$

Let us assume that the forces acting between the various components of the system are *central* in nature, so that \mathbf{f}_{ij} is parallel to $\mathbf{r}_i - \mathbf{r}_j$. In other words, the force exerted on particle j by particle i either points directly toward, or directly away from, particle i , and vice versa. This is a reasonable assumption, because most of the forces that we encounter in the world around us are of this type (e.g., gravity). It follows that if the forces are central in nature then the vector product on the right-hand side of the previous expression is zero. (See Section A.8.) We conclude that

$$\mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_j \times \mathbf{f}_{ji} = \mathbf{0}, \quad (1.86)$$

for all values of i and j . Thus, the double sum on the right-hand side of Equation (1.82) is zero for any kind of central internal force. We are left with

$$\frac{d\mathbf{L}}{dt} = \boldsymbol{\tau}. \quad (1.87)$$

In particular, if the system is isolated, such that it is not subject to a net external torque, so that $\boldsymbol{\tau} = \mathbf{0}$, then the previous equation reduces to

$$\frac{d\mathbf{L}}{dt} = \mathbf{0}. \quad (1.88)$$

In other words, the total angular momentum of an isolated system is a conserved quantity, provided that the different components of the system interact via central forces.

1.5 Invariance Laws

1.5.1 Inertial Reference Frame

Suppose that we have found an inertial frame of reference, and have set up a Cartesian coordinate system in this frame. The motion of particle i in the many-particle system discussed in Section 1.4 is specified by giving its displacement, $\mathbf{r}_i \equiv (x_i, y_i, z_i)$, with respect to the origin of the coordinate system, as a function of time, t . In particular, the linear and angular equations of motion the particle take the respective forms

$$m_i \frac{d\mathbf{v}_i}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{f}_{ij} + \mathbf{F}_i, \quad (1.89)$$

and

$$m_i \frac{d(\mathbf{r}_i \times \mathbf{v}_i)}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_i \times \mathbf{F}_i. \quad (1.90)$$

[See Equations (1.63), (1.79), and (1.81).]

1.5.2 Rotational Invariance

The displacement, \mathbf{r}_i , is, by definition, a vector (i.e., its components, x_i, y_i, z_i , transform under rotation of the coordinate axes in such a manner that the length and direction of \mathbf{r}_i are preserved). (See Section A.5.) Moreover, m_i and t are scalars (i.e., they are invariant under rotation of the coordinate axes). It follows that $\mathbf{v}_i = d\mathbf{r}_i/dt$ and $d\mathbf{v}_i/dt$ are vectors. Furthermore, we have already seen that forces are vectors. (See Section 1.2.3.) Finally, we know that if \mathbf{a} and \mathbf{b} are vectors then so is $\mathbf{a} \times \mathbf{b}$. (See Section A.8.) It follows that every term appearing in the previous two equations transforms as a vector under rotation of the coordinate axes. In other words, the forms of the linear and angular equations of motion, (1.89) and (1.90), respectively, are invariant under rotation of the coordinate axes. Of course, this must be the case because the choice of the orientation of the axes of a Cartesian coordinate system is completely arbitrary, and has no bearing on the motions of bodies in the universe.

1.5.3 Translational Invariance

Suppose that we transform our coordinate system such that the origin shifts from $\mathbf{r} = \mathbf{0}$ to $\mathbf{r} = \mathbf{r}_{\text{shift}}$, where $\mathbf{r}_{\text{shift}}$ is independent of time. It follows that

$$\mathbf{r}_i \rightarrow \mathbf{r}_i - \mathbf{r}_{\text{shift}}, \quad (1.91)$$

$$\mathbf{v}_i \rightarrow \mathbf{v}_i, \quad (1.92)$$

$$\mathbf{f}_{ij} \rightarrow \mathbf{f}_{ij}, \quad (1.93)$$

$$\mathbf{F}_i \rightarrow \mathbf{F}_i. \quad (1.94)$$

The latter two equations follow because forces are obviously not affected by the transformation. It is clear that the linear equation of motion, (1.89), is invariant under the transformation. On the other hand, the angular equation of motion, (1.90), becomes

$$m_i \frac{d(\mathbf{r}_i \times \mathbf{v}_i)}{dt} - m_i \mathbf{r}_{\text{shift}} \times \frac{d\mathbf{v}_i}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_i \times \mathbf{F}_i - \mathbf{r}_{\text{shift}} \times \sum_{j=1, N}^{j \neq i} \mathbf{f}_{ij} - \mathbf{r}_{\text{shift}} \times \mathbf{F}_i. \quad (1.95)$$

However, the vector product of $\mathbf{r}_{\text{shift}}$ with Equation (1.89) yields

$$m_i \mathbf{r}_{\text{shift}} \times \frac{d\mathbf{v}_i}{dt} = \mathbf{r}_{\text{shift}} \times \sum_{j=1, N}^{j \neq i} \mathbf{f}_{ij} + \mathbf{r}_{\text{shift}} \times \mathbf{F}_i, \quad (1.96)$$

The previous two equations can be combined to give

$$m_i \frac{d(\mathbf{r}_i \times \mathbf{v}_i)}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_i \times \mathbf{F}_i. \quad (1.97)$$

Thus, we conclude that the angular equation of motion, (1.90), is also invariant under the transformation. Of course, all of this makes sense because the choice of the origin of a Cartesian coordinate

system is completely arbitrary, and has no bearing on the motions of bodies in the universe. One corollary of the previous analysis is that it does not matter about which point we choose to take moments of momenta and forces to generate angular momenta and torques, respectively, as long as we choose the same point in all cases.

1.5.4 Galilean Invariance

Consider a second frame of reference moving with some arbitrary constant velocity \mathbf{u} with respect to our original inertial reference frame. We can assume, without loss of generality, that the origins of the two coordinate systems coincide at time $t = 0$. (Note that there is a tacit assumption that clocks run at the same rate in both reference frames. This is the case provided that the relative speed of the two frames is much smaller than the speed of light in vacuum. See Section 3.2.3.) It follows that

$$\mathbf{r}_i \rightarrow \mathbf{r}_i - \mathbf{u}t, \quad (1.98)$$

$$\mathbf{v}_i \rightarrow \mathbf{v}_i - \mathbf{u}, \quad (1.99)$$

$$\frac{d\mathbf{v}_i}{dt} \rightarrow \frac{d\mathbf{v}_i}{dt}, \quad (1.100)$$

$$\mathbf{f}_{ij} \rightarrow \mathbf{f}_{ij}, \quad (1.101)$$

$$\mathbf{F}_i \rightarrow \mathbf{F}_i. \quad (1.102)$$

Here, we are assuming that the forces are the same in both reference frames. It is clear that the linear equation of motion, (1.89), takes the same form in the second reference frame. On the other hand, the angular equation of motion, (1.90), becomes

$$m_i \frac{d(\mathbf{r}_i \times \mathbf{v}_i)}{dt} - m_i t \mathbf{u} \times \frac{d\mathbf{v}_i}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_i \times \mathbf{F}_i - t \mathbf{u} \times \sum_{j=1, N}^{j \neq i} \mathbf{f}_{ij} - t \mathbf{u} \times \mathbf{F}_i. \quad (1.103)$$

However, the vector product of \mathbf{u} with Equation (1.89) yields

$$m_i \mathbf{u} \times \frac{d\mathbf{v}_i}{dt} = \mathbf{u} \times \sum_{j=1, N}^{j \neq i} \mathbf{f}_{ij} + \mathbf{u} \times \mathbf{F}_i, \quad (1.104)$$

The previous two equations can be combined to give

$$m_i \frac{d(\mathbf{r}_i \times \mathbf{v}_i)}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_i \times \mathbf{F}_i. \quad (1.105)$$

Thus the angular equation of motion, (1.90), also takes the same form in the second reference frame. It is clear that the second frame of reference is also an inertial reference frame (i.e., a frame in which Newton's laws of motion are valid). In particular, if all of the forces are zero then the

system's constituent particle do not accelerate in either the first or the second reference frame, in accordance with Newton's first law of motion. Given that the constant velocity \mathbf{u} is arbitrary, we conclude that there are an infinite number of different inertial frames of reference all moving at constant velocities with respect to one another, and that Newton's laws of motion are equally valid in each frame. In particular, if the total energy of the system is conserved in one inertial reference frame then it is conserved in all inertial reference frames. Likewise, if the total momentum of the system is conserved in one inertial reference frame then it is conserved in all inertial reference frames. Finally, if the total angular momentum of the system is conserved in one inertial reference frame then it is conserved in all inertial reference frames.

Consider the special case in which $\mathbf{u} = u \mathbf{e}_x$. Let $\mathbf{r} = (x, y, z)$ be a general displacement in the first reference frame, and let $\mathbf{r}' = (x', y', z')$ be the corresponding displacement in the second frame. It follows from Equation (1.98) that

$$x' = x - u t, \quad (1.106)$$

$$y' = y, \quad (1.107)$$

$$z' = z. \quad (1.108)$$

This coordinate transformation was implied in the researches of Galileo Galilei, and is known as a *Galilean transformation* in his honor. Hence, the fact that Newton's laws of motion take the same form in all inertial reference frames is known as *Galilean invariance*.

Suppose that the second frame of reference accelerates with respect to the first. In other words, suppose that $\mathbf{u} = \mathbf{u}(t)$. It follows that

$$\mathbf{r}_i \rightarrow \mathbf{r}_i - \int_0^t \mathbf{u}(t') dt', \quad (1.109)$$

$$\mathbf{v}_i \rightarrow \mathbf{v}_i - \mathbf{u}, \quad (1.110)$$

$$\frac{d\mathbf{v}_i}{dt} \rightarrow \frac{d\mathbf{v}_i}{dt} - \frac{d\mathbf{u}}{dt}, \quad (1.111)$$

$$\mathbf{f}_{ij} \rightarrow \mathbf{f}_{ij}, \quad (1.112)$$

$$\mathbf{F}_i \rightarrow \mathbf{F}_i. \quad (1.113)$$

It is easily seen that Equations (1.89) and (1.90) transform to give

$$m_i \frac{d\mathbf{v}_i}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{f}_{ij} + \mathbf{F}_i + m_i \frac{d\mathbf{u}}{dt}, \quad (1.114)$$

$$m_i \frac{d(\mathbf{r}_i \times \mathbf{v}_i)}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{r}_i \times \mathbf{f}_{ij} + \mathbf{r}_i \times \mathbf{F}_i + m_i \mathbf{r}_i \times \frac{d\mathbf{u}}{dt}, \quad (1.115)$$

respectively. Note the appearance of the so-called *fictitious force*, $m_i d\mathbf{u}/dt$, and the so-called *fictitious torque*, $m_i \mathbf{r}_i \times d\mathbf{u}/dt$, on the right-hand sides of the previous two equations. It is clear that

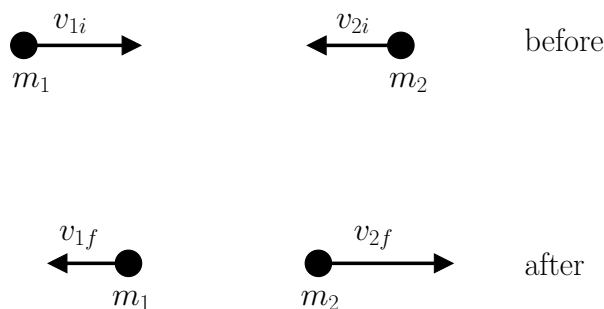


Figure 1.4: A one-dimension collision in the laboratory frame.

the second frame of reference is not an inertial frame. (For instance, if all of the real forces were zero then the system's constituent particles would still accelerate, which is not in accordance with Newton's first law of motion.) Hence, we conclude that any frame of reference that accelerates with respect to a given inertial reference frame is non-inertial.

1.6 Two-Particle Collisions

1.6.1 One-Dimensional Collisions

Consider two particles of mass m_1 and m_2 , respectively, that are free to move in one dimension. Suppose that these two particles collide. Suppose, further, that both particles are subject to zero net force when they are not in contact with one another. Finally, let us assume that we are observing the collision in a convenient inertial reference frame known as the *laboratory frame*. This situation is illustrated in Figure 1.4.

Both before and after the collision, the two particles move with constant velocity, in accordance with Newton's first law of motion. Let v_{1i} and v_{2i} be the velocities of the first and second particles, respectively, before the collision. Here, velocities to the right in Figure 1.4 are positive. Likewise, let v_{1f} and v_{2f} be the velocities of the first and second particles, respectively, after the collision. During the collision itself, the first particle exerts a large transitory force, f_{21} , on the second, whereas the second particle exerts an equal and opposite force, $f_{12} = -f_{21}$, on the first. In fact, we can model the collision as equal and opposite impulses given to the two particles at the instant in time when they come together. (See Section 1.3.1.)

We are clearly considering a system in which there is zero net external force (because the forces associated with the collision are internal in nature). Hence, the total (linear) momentum of the system is a conserved quantity. (See Section 1.4.4.) Equating the total momenta before and after the collision, we obtain

$$m_1 v_{1i} + m_2 v_{2i} = m_1 v_{1f} + m_2 v_{2f}. \quad (1.116)$$

This equation is valid for any one-dimensional collision, irrespective its nature.

Suppose that the collision is *elastic*, which means that there is no associated loss of kinetic energy. Equating the net kinetic energies before and after the collision, we obtain

$$\frac{1}{2} m_1 v_{1i}^2 + \frac{1}{2} m_2 v_{2i}^2 = \frac{1}{2} m_1 v_{1f}^2 + \frac{1}{2} m_2 v_{2f}^2. \quad (1.117)$$

(See Section 1.3.2.) It follows that

$$m_1 (v_{1f}^2 - v_{1i}^2) = -m_2 (v_{2f}^2 - v_{2i}^2), \quad (1.118)$$

or

$$m_1 (v_{1f} - v_{1i})(v_{1f} + v_{1i}) = -m_2 (v_{2f} - v_{2i})(v_{2f} + v_{2i}). \quad (1.119)$$

However, Equation (1.116) yields

$$m_1 (v_{1f} - v_{1i}) = -m_2 (v_{2f} - v_{2i}). \quad (1.120)$$

The previous two equations can be combined to give

$$v_{1f} + v_{1i} = v_{2f} + v_{2i}, \quad (1.121)$$

or

$$(v_{2f} - v_{1f}) = -(v_{2i} - v_{1i}). \quad (1.122)$$

Thus, we conclude that an elastic collision causes the relative velocity of the two particles to reverse direction, while keeping the same magnitude.

Suppose that we transform to a frame of reference that co-moves with the center of mass of the system. The motion of a multi-particle system often looks particularly simple when viewed in such a frame. Because the system is subject to zero net external force, the velocity of the center of mass is invariant [see Equations (1.72)], and is given by

$$V = \frac{m_1 v_{1i} + m_2 v_{2i}}{m_1 + m_2} = \frac{m_1 v_{1f} + m_2 v_{2f}}{m_1 + m_2}. \quad (1.123)$$

[See Equation (1.68).] A particle that possesses a velocity v in the laboratory frame possesses a velocity $v' = v - V$ in the so-called *center-of-mass frame*. It is easily demonstrated that

$$v'_{1i} = -\frac{m_2}{m_1 + m_2} (v_{2i} - v_{1i}), \quad (1.124)$$

$$v'_{2i} = +\frac{m_1}{m_1 + m_2} (v_{2i} - v_{1i}), \quad (1.125)$$

$$v'_{1f} = -\frac{m_2}{m_1 + m_2} (v_{2f} - v_{1f}), \quad (1.126)$$

$$v'_{2f} = +\frac{m_1}{m_1 + m_2} (v_{2f} - v_{1f}). \quad (1.127)$$

Note, incidentally, that the center-of-mass frame is obviously inertial (because it is moving at a constant velocity with respect to the inertial laboratory frame).

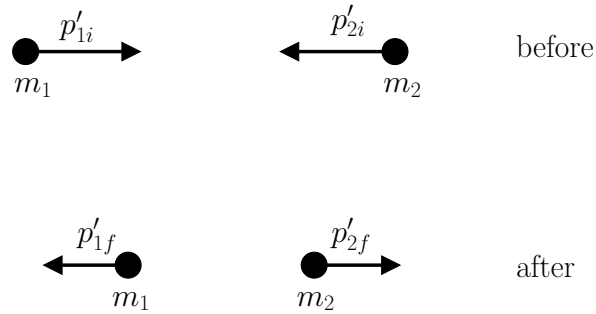


Figure 1.5: A one-dimension collision in the center-of-mass frame.

The previous four equations yield

$$-p'_{1i} = p'_{2i} = \mu (v_{2i} - v_{1i}), \quad (1.128)$$

$$-p'_{1f} = p'_{2f} = \mu (v_{2f} - v_{1f}), \quad (1.129)$$

where $\mu = m_1 m_2 / (m_1 + m_2)$ is the so-called *reduced mass* (see Section 1.10.7), and $p'_{1i} = m_1 v'_{1i}$ is the initial momentum of the first particle in the center-of-mass frame, et cetera. In other words, when viewed in the center-of-mass frame, the two particles approach one another with equal and opposite momenta before the collision, and diverge from one another with equal and opposite momenta after the collision. See Figure 1.5. Thus, the center-of-mass momentum conservation equation,

$$p'_{1i} + p'_{2i} = p'_{1f} + p'_{2f}, \quad (1.130)$$

is trivially satisfied, because both the left- and right-hand sides are zero. Incidentally, this result is valid for both elastic and inelastic collisions.

Equations (1.122), (1.128), and (1.129) can be combined to give

$$p'_{1f} = -p'_{1i}, \quad (1.131)$$

$$p'_{2f} = -p'_{2i}. \quad (1.132)$$

In other words, in the center-of-mass frame, an elastic collision causes the equal and opposite momenta of the two particles to both reverse direction, but keep the same magnitude. The previous two expressions imply that

$$v'_{1f} = -v'_{1i}, \quad (1.133)$$

$$v'_{2f} = -v'_{2i}. \quad (1.134)$$

In other words, in the center-of-mass frame, an elastic collision also causes the velocity of each particle to reverse direction, but keep the same magnitude. Thus, the total kinetic energy of the system is obviously a conserved quantity in the center-of-mass frame.

Equations (1.124) and (1.125) can be combined with the previous two equations to give

$$v'_{1f} = \frac{m_2}{m_1 + m_2} (v_{2i} - v_{1i}), \quad (1.135)$$

$$v'_{2f} = -\frac{m_1}{m_1 + m_2} (v_{2i} - v_{1i}). \quad (1.136)$$

However, $v_{1f} = v'_{1f} + V$ and $v_{2f} = v'_{2f} + V$, which allows us to express the velocities of the two particles after the collision in the laboratory frame in terms of the corresponding velocities before the collision:

$$v_{1f} = \left(\frac{m_1 - m_2}{m_1 + m_2} \right) v_{1i} + \left(\frac{2m_2}{m_1 + m_2} \right) v_{2i}, \quad (1.137)$$

$$v_{2f} = \left(\frac{2m_1}{m_1 + m_2} \right) v_{1i} - \left(\frac{m_1 - m_2}{m_1 + m_2} \right) v_{2i}. \quad (1.138)$$

Let us, now, consider some special cases. Suppose that two equal-mass particles collide elastically. If $m_1 = m_2$ then Equations (1.137) and (1.138) yield

$$v_{1f} = v_{2i}, \quad (1.139)$$

$$v_{2f} = v_{1i}. \quad (1.140)$$

In other words, the two particles simply exchange velocities when they collide. For instance, if the second particle is stationary and the first particle strikes it head-on with velocity v then the first particle is brought to a halt whereas the second particle moves off with velocity v . It is possible to reproduce this effect in snooker or pool by striking the cue ball with great force in such a manner that it slides, rather than rolls, over the table; in this case, when the cue ball strikes another ball head-on it comes to a complete halt, and the other ball is propelled forward very rapidly. Incidentally, it is necessary to prevent the cue ball from rolling, because rolling motion is not taken into account in our analysis, and actually changes the answer.

Suppose that the second particle is much more massive than the first (i.e., $m_2 \gg m_1$), and is initially at rest (i.e., $v_{2i} = 0$). In this case, Equations (1.137) and (1.138) yield

$$v_{1f} \simeq -v_{1i}, \quad (1.141)$$

$$v_{2f} \simeq 0. \quad (1.142)$$

In other words, the velocity of the light particle is effectively reversed during the collision, whereas the massive particle remains approximately at rest. Indeed, this is the sort of behavior we expect when an object collides elastically with an immovable obstacle; for instance, when an elastic ball bounces off a brick wall.

Suppose, finally, that the second particle is much lighter than the first (i.e., $m_2 \ll m_1$), and is initially at rest (i.e., $v_{2i} = 0$). In this case, Equations (1.137) and (1.138) yield

$$v_{1f} \simeq v_{1i}, \quad (1.143)$$

$$v_{2f} \simeq 2v_{1i}. \quad (1.144)$$

In other words, the motion of the massive particle is essentially unaffected by the collision, whereas the light particle ends up moving twice as fast as the massive one.

1.6.2 Totally Inelastic Collisions

In a *totally inelastic* collision, the two particles stick together after colliding, so that they end up moving with the same final velocity, $v_f = v_{1f} = v_{2f}$. In this case,

$$v_f = \frac{m_1 v_{1i} + m_2 v_{2i}}{m_1 + m_2} = V. \quad (1.145)$$

In other words, the common final velocity of the two particles is equal to the center-of-mass velocity of the system. This is hardly a surprising result. We have already seen that in the center-of-mass frame the two particles must diverge with equal and opposite momenta after the collision. However, in a totally inelastic collision these two momenta must also be equal (because the two objects stick together). The only way in which this is possible is if the two particles remain stationary in the center-of-mass frame after the collision. Hence, after the collision, the two particles move with the center-of-mass velocity in the laboratory frame.

Suppose that the second object is initially at rest (i.e., $v_{2i} = 0$) in the laboratory frame. In this special case, the common final velocity of the two objects is

$$v_f = \frac{m_1}{m_1 + m_2} v_{1i}. \quad (1.146)$$

Note that the first object is slowed down by the collision. The fractional loss in kinetic energy of the system due to the collision is given by

$$f = \frac{m_2}{m_1 + m_2}. \quad (1.147)$$

The loss in kinetic energy is small if the (initially) stationary object is much lighter than the moving object (i.e., if $m_2 \ll m_1$), and almost 100% if the moving object is much lighter than the stationary one (i.e., if $m_2 \gg m_1$). Of course, the lost kinetic energy of the system is converted into some other form of energy; for instance, heat energy.

1.6.3 Two-Dimensional Collisions

Suppose that an object of mass m_1 , moving with initial velocity \mathbf{v}_{1i} , strikes a second object, of mass m_2 , that is initially at rest. Suppose, further, that the collision is not head-on, so that after the collision the first object moves off at an angle θ_1 to its initial direction of motion, whereas the second object recoils at an angle θ_2 to this direction. Let the final velocities of the two objects be \mathbf{v}_{1f} and \mathbf{v}_{2f} , respectively. See Figure 1.6.

We are again considering a system in which there is zero net external force (because the forces associated with the collision are internal in nature). It follows that the total momentum of the system is a conserved quantity. However, unlike before, we must now treat momentum as a vector quantity, because we are no longer dealing with one-dimensional motion. Momentum conservation implies that

$$m_2 \mathbf{v}_{1i} = m_1 \mathbf{v}_{1f} + m_2 \mathbf{v}_{2f}. \quad (1.148)$$

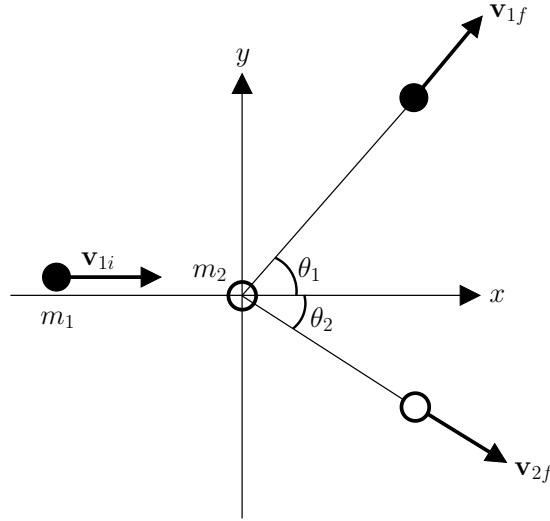


Figure 1.6: A two-dimensional collision in the laboratory frame.

As before, it is convenient to transform to a frame of reference that co-moves with the center of mass of the system. The invariant velocity of the center of mass is given by

$$\mathbf{V} = \frac{m_1 \mathbf{v}_{1i}}{m_1 + m_2} = \frac{m_1 \mathbf{v}_{1f} + m_2 \mathbf{v}_{2f}}{m_1 + m_2}. \quad (1.149)$$

An object that possesses a velocity \mathbf{v} in the laboratory frame possesses a velocity $\mathbf{v}' = \mathbf{v} - \mathbf{V}$ in the center-of-mass frame. Hence, it follows that

$$\mathbf{v}'_{1i} = \left(\frac{m_2}{m_1 + m_2} \right) \mathbf{v}_{1i}, \quad (1.150)$$

$$\mathbf{v}'_{2i} = - \left(\frac{m_1}{m_1 + m_2} \right) \mathbf{v}_{1i}, \quad (1.151)$$

$$\mathbf{v}'_{1f} = - \left(\frac{m_2}{m_1 + m_2} \right) (\mathbf{v}_{2f} - \mathbf{v}_{1f}), \quad (1.152)$$

$$\mathbf{v}'_{2f} = \left(\frac{m_1}{m_1 + m_2} \right) (\mathbf{v}_{2f} - \mathbf{v}_{1f}). \quad (1.153)$$

Furthermore, the momenta in the center-of-mass frame take the form

$$-\mathbf{p}'_{1i} = \mathbf{p}'_{2i} = -\mu \mathbf{v}_{1i}, \quad (1.154)$$

$$-\mathbf{p}'_{1f} = \mathbf{p}'_{2f} = \mu (\mathbf{v}_{2f} - \mathbf{v}_{1f}), \quad (1.155)$$

where $\mu = m_1 m_2 / (m_1 + m_2)$. (Of course, $\mathbf{p}'_{1i} = m_1 \mathbf{v}'_{1i}$, et cetera.) As before, in the center-of-mass frame, the two objects approach one another with equal and opposite momenta before the collision, and diverge from one another with equal and opposite momenta after the collision. Let θ be the

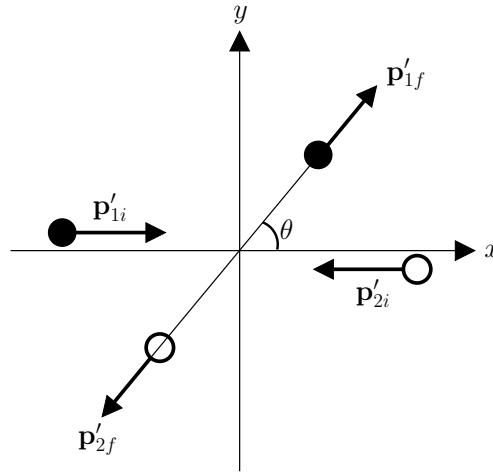


Figure 1.7: A two-dimensional collision in the center-of-mass frame.

direction subtended between the final and initial momenta of each object in the center-of-mass frame. See Figure 1.7. It follows that in the x - y coordinate system shown in the figure,

$$\mathbf{p}'_{1i} = p'_{1i} (1, 0), \quad (1.156)$$

$$\mathbf{p}'_{2i} = -p'_{1i} (1, 0), \quad (1.157)$$

$$\mathbf{p}'_{1f} = p'_{1f} (\cos \theta, \sin \theta), \quad (1.158)$$

$$\mathbf{p}'_{2f} = -p'_{1f} (\cos \theta, \sin \theta), \quad (1.159)$$

where $p'_{1i} = |\mathbf{p}'_{1i}|$, et cetera. Finally, if the collision is elastic then the kinetic energy before the collision must equal that after the collision (in both the laboratory and the center-of-mass frames). It follows that

$$p'_{1f} = p'_{1i}. \quad (1.160)$$

Hence,

$$\mathbf{v}'_{1f} \equiv \frac{\mathbf{p}'_{1f}}{m_1} = \left(\frac{m_2 v_{1i}}{m_1 + m_2} \right) (\cos \theta, \sin \theta), \quad (1.161)$$

$$\mathbf{v}'_{2f} \equiv \frac{\mathbf{p}'_{2f}}{m_2} = - \left(\frac{m_1 v_{1i}}{m_1 + m_2} \right) (\cos \theta, \sin \theta), \quad (1.162)$$

It can be seen from the previous two equations that an elastic two-dimensional collision is fully characterized once the initial velocity, v_{1i} , and the scattering angle, θ , are specified. In general, we would expect θ to be able to take all values in the range 0 to π . In fact, a head-on collision corresponds to $\theta = \pi$, whereas a glancing collision corresponds to $\theta \ll 1$.

Now, $\mathbf{v} = \mathbf{v}' + \mathbf{V}$, where

$$\mathbf{V} = \left(\frac{m_1 v_{1i}}{m_1 + m_2} \right) (1, 0). \quad (1.163)$$

It follows that, in the x - y coordinate system shown in Figure 1.6, the laboratory-frame velocities of the two objects after the collision are

$$\mathbf{v}_{1f} = \left(\frac{v_{1i}}{m_1 + m_2} \right) (m_1 + m_2 \cos \theta, m_2 \sin \theta), \quad (1.164)$$

$$\mathbf{v}_{2f} = \left(\frac{m_1 v_{1i}}{m_1 + m_2} \right) (1 - \cos \theta, -\sin \theta). \quad (1.165)$$

Hence, according to Figure 1.6,

$$\tan \theta_1 = \frac{\sin \theta}{\cos \theta + m_1/m_2}, \quad (1.166)$$

$$\tan \theta_2 = \frac{\sin \theta}{1 - \cos \theta} = \tan \left(\frac{\pi}{2} - \frac{\theta}{2} \right). \quad (1.167)$$

The last equation implies that

$$\theta_2 = \frac{\pi}{2} - \frac{\theta}{2}. \quad (1.168)$$

Differentiating Equation (1.166) with respect to θ , we obtain

$$\frac{d \tan \theta_1}{d\theta} = \frac{1 + (m_1/m_2) \cos \theta}{(\cos \theta + m_1/m_2)^2}. \quad (1.169)$$

Thus, $\tan \theta_1$ attains an extreme value, which can be shown to correspond to a maximum possible value of θ_1 , when the numerator of the previous expression is zero; that is, when

$$\cos \theta = -\frac{m_2}{m_1}. \quad (1.170)$$

Note that it is only possible to solve the previous equation when $m_1 > m_2$. If this is the case then Equation (1.166) yields

$$\tan \theta_{1 \max} = \frac{m_2/m_1}{\sqrt{1 - (m_2/m_1)^2}}, \quad (1.171)$$

which reduces to

$$\theta_{1 \max} = \sin^{-1} \left(\frac{m_2}{m_1} \right). \quad (1.172)$$

Hence, we conclude that when $m_1 > m_2$ there is a maximum possible value of the scattering angle, θ_1 , in the laboratory frame. This maximum value is always less than $\pi/2$, which implies that there is no backward scattering (i.e., $\theta_1 > \pi/2$) at all when $m_1 > m_2$. For the special case when $m_1 = m_2$, the maximum scattering angle is $\pi/2$. However, for $m_1 < m_2$ there is no maximum value, and the scattering angle in the laboratory frame can thus range all the way to π .

Suppose that the two particles have equal masses, so that $m_1 = m_2$. In this case, Equation (1.166) yields

$$\tan \theta_1 = \frac{\sin \theta}{\cos \theta + 1} = \tan \left(\frac{\theta}{2} \right). \quad (1.173)$$

Hence,

$$\theta_1 = \frac{\theta}{2}. \quad (1.174)$$

In other words, the scattering angle of the first particle in the laboratory frame is half of the scattering angle in the center-of-mass frame. The previous equation can be combined with Equation (1.168) to give

$$\theta_1 + \theta_2 = \frac{\pi}{2}. \quad (1.175)$$

Thus, in the laboratory frame, the two particles move off at right-angles to one another after the collision. It is possible to reproduce this effect in snooker or pool by striking the cue ball with great force in such a manner that it slides, rather than rolls, over the table; in this case, when the cue ball strikes another ball obliquely then the two balls move off at right-angles to one another. Incidentally, it is necessary to prevent the cue ball from rolling, because rolling motion is not taken into account in our analysis, and actually changes the answer. Finally, it is easily demonstrated that the fractions of the initial kinetic energy carried off by the two particles after the collision are

$$\frac{E_{1f}}{E_{1i}} = \cos^2 \theta_1, \quad (1.176)$$

$$\frac{E_{2f}}{E_{1i}} = \sin^2 \theta_1. \quad (1.177)$$

1.7 Rigid Body Rotation

1.7.1 Fundamental Equations

We can think of a rigid body as a collection of a large number of small mass elements that all maintain a fixed spatial relationship with respect to one another. Let there be N elements, and let the i th element have mass m_i , instantaneous displacement \mathbf{r}_i , and instantaneous velocity \mathbf{v}_i . The equation of motion of the i th element is written

$$m_i \frac{d\mathbf{v}_i}{dt} = \sum_{j=1, N}^{j \neq i} \mathbf{f}_{ij} + \mathbf{F}_i. \quad (1.178)$$

(See Section 1.4.1.) Here, \mathbf{f}_{ij} is the internal force exerted on the i th element by the j th element, and \mathbf{F}_i the external force acting on the i th element. The internal forces, \mathbf{f}_{ij} , represent the stresses that develop within the body in order to ensure that its various elements maintain a constant spatial relationship with respect to one another. Of course, $\mathbf{f}_{ij} = -\mathbf{f}_{ji}$, by Newton's third law. The external forces represent forces that originate outside the body.

Repeating the analysis of Section 1.4.2, we can sum Equation (1.178) over all mass elements to obtain

$$M \frac{d^2 \mathbf{R}}{dt^2} = \mathbf{F}. \quad (1.179)$$

Here, $M = \sum_{i=1, N} m_i$ is the total mass, \mathbf{R} the displacement of the center of mass [see Equation (1.68)], and $\mathbf{F} = \sum_{i=1, N} \mathbf{F}_i$ the total external force. It can be seen that the center of mass

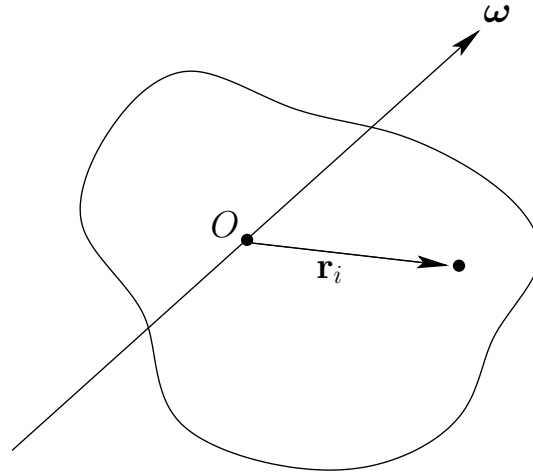


Figure 1.8: A rigid rotating body.

of a rigid body moves under the action of the external forces like a point particle whose mass is identical to that of the body.

Repeating the analysis of Section 1.4.5, we can sum $\mathbf{r}_i \times$ Equation (1.178) over all mass elements to obtain

$$\frac{d\mathbf{L}}{dt} = \boldsymbol{\tau}. \quad (1.180)$$

Here, $\mathbf{L} = \sum_{i=1,N} m_i \mathbf{r}_i \times \mathbf{v}_i$ is the total angular momentum of the body (about the origin), and $\boldsymbol{\tau} = \sum_{i=1,N} \mathbf{r}_i \times \mathbf{F}_i$ the total external torque (about the origin). Note that the previous equation is only valid if the internal forces are central in nature. However, this is not a particularly onerous constraint. Equation (1.180) describes how the angular momentum of a rigid body evolves in time under the action of the external torques.

In the following, we shall only consider the rotational motion of rigid bodies, because their translational motion is similar to that of point particles [see Equation (1.179)], and, therefore, fairly straightforward in nature.

1.7.2 Moment of Inertia Tensor

Consider a rigid body rotating with fixed angular velocity $\boldsymbol{\omega}$ about an axis that passes through the origin. See Figure 1.8. Let \mathbf{r}_i be the displacement of the i th mass element, whose mass is m_i . We expect this displacement to precess about the axis of rotation (which is parallel to $\boldsymbol{\omega}$) with angular velocity $\boldsymbol{\omega}$. It, therefore, follows from Equation (A.52) that

$$\mathbf{v}_i = \boldsymbol{\omega} \times \mathbf{r}_i. \quad (1.181)$$

The total angular momentum of the body (about the origin) is written

$$\mathbf{L} = \sum_{i=1,N} m_i \mathbf{r}_i \times \mathbf{v}_i = \sum_{i=1,N} m_i \mathbf{r}_i \times (\boldsymbol{\omega} \times \mathbf{r}_i) = \sum_{i=1,N} m_i [r_i^2 \boldsymbol{\omega} - (\mathbf{r}_i \cdot \boldsymbol{\omega}) \mathbf{r}_i], \quad (1.182)$$

where use has been made of Equation (1.181), and some standard vector identities. (See Section A.11.) The previous formula can be written as a matrix equation of the form

$$\begin{pmatrix} L_x \\ L_y \\ L_z \end{pmatrix} = \begin{pmatrix} I_{xx} & I_{xy} & I_{xz} \\ I_{yx} & I_{yy} & I_{yz} \\ I_{zx} & I_{zy} & I_{zz} \end{pmatrix} \begin{pmatrix} \omega_x \\ \omega_y \\ \omega_z \end{pmatrix}, \quad (1.183)$$

where

$$I_{xx} = \sum_{i=1,N} (y_i^2 + z_i^2) m_i, \quad (1.184)$$

$$I_{yy} = \sum_{i=1,N} (x_i^2 + z_i^2) m_i, \quad (1.185)$$

$$I_{zz} = \sum_{i=1,N} (x_i^2 + y_i^2) m_i, \quad (1.186)$$

$$I_{xy} = I_{yx} = - \sum_{i=1,N} x_i y_i m_i \quad (1.187)$$

$$I_{yz} = I_{zy} = - \sum_{i=1,N} y_i z_i m_i, \quad (1.188)$$

$$I_{xz} = I_{zx} = - \sum_{i=1,N} x_i z_i m_i. \quad (1.189)$$

Here, I_{xx} is called the *moment of inertia* about the x -axis, I_{yy} the moment of inertia about the y -axis, I_{xy} the x - y *product of inertia*, I_{yz} the y - z product of inertia, et cetera. The matrix of the I_{ij} values is known as the *moment of inertia tensor*.

Suppose that our body is rotationally symmetric about the z -axis. In this case, it is easily seen that the products of inertia are all zero. Moreover, $I_{xx} = I_{yy} = I_{\perp}$. Let us write $I_{zz} = I_{\parallel}$. Note that, in general, $I_{\parallel} \neq I_{\perp}$ (unless the body is spherically symmetric). Thus, Equation (1.183) simplifies to give

$$\mathbf{L} = I_{\perp} \omega_x \mathbf{e}_x + I_{\perp} \omega_y \mathbf{e}_y + I_{\parallel} \omega_z \mathbf{e}_z. \quad (1.190)$$

The angular momentum vector, \mathbf{L} , obtained from the previous equation, does not necessarily point in the same direction as the angular velocity vector, $\boldsymbol{\omega}$ (because $I_{\parallel} \neq I_{\perp}$). In other words, \mathbf{L} is generally not parallel to $\boldsymbol{\omega}$. However, if the body rotates about \mathbf{e}_z or any axis in the x - y plane then \mathbf{L} is parallel to $\boldsymbol{\omega}$. These special axes of rotation are called *principal axes of rotation*, and the associated moments of inertia, I_{\parallel} and I_{\perp} , respectively, are called *principal moments of inertia*.

It can be demonstrated that any rigid body (not just an axisymmetric one) has three mutually perpendicular principal axes of rotation. Furthermore, if a body is rotating about one of its principal axes of rotation then

$$\mathbf{L} = I \boldsymbol{\omega}, \quad (1.191)$$

where I is the associated principal moment of inertia. More generally, assuming that the Cartesian axes are parallel to the principal moments of inertia, we can write

$$\mathbf{L} = I_{xx} \omega_x \mathbf{e}_x + I_{yy} \omega_y \mathbf{e}_y + I_{zz} \omega_z \mathbf{e}_z, \quad (1.192)$$

where I_{xx} , I_{yy} , and I_{zz} are the three principal moments of inertia.

1.7.3 Rotational Kinetic Energy

The instantaneous rotational kinetic energy of a rotating rigid body is written

$$K = \frac{1}{2} \sum_{i=1,N} m_i \mathbf{v}_i \cdot \mathbf{v}_i. \quad (1.193)$$

Making use of Equation (1.181), and some vector identities (see Section A.10), the kinetic energy takes the form

$$K = \frac{1}{2} \sum_{i=1,N} m_i (\boldsymbol{\omega} \times \mathbf{r}_i) \cdot (\boldsymbol{\omega} \times \mathbf{r}_i) = \frac{1}{2} \boldsymbol{\omega} \cdot \sum_{i=1,N} m_i \mathbf{r}_i \times \boldsymbol{\omega} \times \mathbf{r}_i. \quad (1.194)$$

Hence, it follows from Equation (1.182) that

$$K = \frac{1}{2} \boldsymbol{\omega} \cdot \mathbf{L}. \quad (1.195)$$

For the special case of an axisymmetric body, making use of Equation (1.190), we obtain

$$K = \frac{1}{2} I_{\perp} (\omega_x^2 + \omega_y^2) + \frac{1}{2} I_{\parallel} \omega_z^2. \quad (1.196)$$

For the special case of a body rotating about a principal axis of rotation,

$$K = \frac{1}{2} I \omega^2, \quad (1.197)$$

where I is the associated principal moment of inertia. More generally,

$$K = \frac{1}{2} I_{xx} \omega_x^2 + \frac{1}{2} I_{yy} \omega_y^2 + \frac{1}{2} I_{zz} \omega_z^2 = \frac{L_x^2}{2I_{xx}} + \frac{L_y^2}{2I_{yy}} + \frac{L_z^2}{2I_{zz}}, \quad (1.198)$$

assuming that the Cartesian axes are parallel to the principal axes of rotation.

1.7.4 Power

It follows from Equation (1.193) that

$$\frac{dK}{dt} = \sum_{i=1,N} m_i \frac{d\mathbf{v}_i}{dt} \cdot \mathbf{v}_i. \quad (1.199)$$

Making use of Equation (1.181), we obtain

$$\frac{dK}{dt} = \sum_{i=1,N} m_i \frac{d\mathbf{v}_i}{dt} \cdot \boldsymbol{\omega} \times \mathbf{r}_i = \boldsymbol{\omega} \cdot \sum_{i=1,N} m_i \mathbf{r}_i \times \frac{d\mathbf{v}_i}{dt}. \quad (1.200)$$

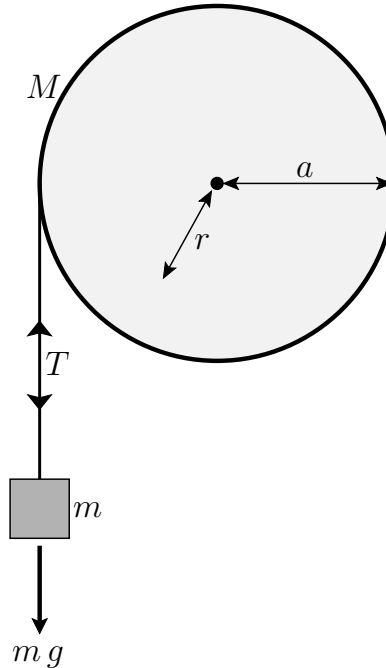


Figure 1.9: A flywheel.

However, according to Equations (1.180) and (1.182),

$$\frac{d\mathbf{L}}{dt} = \frac{d}{dt} \sum_{i=1,N} m_i \mathbf{r}_i \times \mathbf{v}_i = \sum_{i=1,N} m_i \mathbf{r}_i \times \frac{d\mathbf{v}_i}{dt} = \boldsymbol{\tau}. \quad (1.201)$$

Hence, we obtain

$$\frac{dK}{dt} = \boldsymbol{\tau} \cdot \boldsymbol{\omega}. \quad (1.202)$$

The right-hand side of the previous equation specified the work per unit time done on the system by the external torque.

1.7.5 Uniform Flywheel

Consider a uniform flywheel of mass M and radius a . Suppose that the flywheel rotates about a horizontal axis that passes through its center, and is perpendicular to the plane of the flywheel. Let the flywheel have a light inextensible cord wrapped around its circumference to one end of which is attached a mass m that dangles below the flywheel. Let T be the tension in the cord. See Figure 1.9. Let us determine the angular acceleration of the flywheel.

The flywheel is rotationally symmetric about its axis of rotation, which implies that this axis is a principal axis of rotation. Suppose that the rotation axis corresponds to the z -axis (this axis is directed out of the paper in Figure 1.9). It follows that

$$\boldsymbol{\omega} = \omega \mathbf{e}_z, \quad (1.203)$$

$$\mathbf{L} = I \omega \mathbf{e}_z, \quad (1.204)$$

where I is the associated principal moment of inertia, and use has been made of Equation (1.191). Now, it is clear that the tension in the cord exerts a torque

$$\boldsymbol{\tau} = T a \mathbf{e}_z \quad (1.205)$$

on the flywheel. Hence, the rotational equation of motion of the flywheel, (1.180), yields

$$I \frac{d\omega}{dt} = \tau = T a. \quad (1.206)$$

Consider the equation of motion of the dangling mass. We can write

$$m \frac{dv}{dt} = m g - T, \quad (1.207)$$

where v is the mass's downward velocity, and g is the acceleration due to gravity. Because the cord is inextensible, it follows that v is also the downward velocity of the cord. Assuming that the cord unwraps from the flywheel without slipping, its downward velocity must match the tangential velocity of the flywheel's outer rim. In other words,

$$v = a \omega, \quad (1.208)$$

which implies that

$$\frac{dv}{dt} = a \frac{d\omega}{dt}. \quad (1.209)$$

Equations (1.206), (1.207), and (1.209) yield

$$T = \left(\frac{I}{I + m a^2} \right) m g, \quad (1.210)$$

$$\frac{dv}{dt} = \left(\frac{m a^2}{I + m a^2} \right) g, \quad (1.211)$$

$$\frac{d\omega}{dt} = \left(\frac{m a^2}{I + m a^2} \right) \frac{g}{a}. \quad (1.212)$$

It remains to calculate the principal moment of inertia, I , of the flywheel about its rotation axis. The moment of inertia is written

$$I = \sum_i m_i r_i^2 \quad (1.213)$$

for all of the mass elements that make up then flywheel. Here, r represents radial distance from the axis of rotation. See Figure 1.9. We can assume that the flywheel has a constant mass per unit area (in the plane perpendicular to the rotation axis), ρ . Consider the contribution of an annulus of inner radius r and outer radius $r + dr$ to the moment of inertia. It is clear that the area of the annulus is $2\pi r dr$, Thus, its mass is $dm = 2\pi r dr \rho$. Hence,

$$dI = dm r^2 = 2\pi \rho r^3 dr, \quad (1.214)$$

which implies that

$$I = 2\pi\rho \int_0^a r^3 dr = \frac{\pi\rho a^4}{2}. \quad (1.215)$$

However, the total mass of the flywheel is

$$M = \int dm = 2\pi\rho \int_0^a r dr = \pi\rho a^2. \quad (1.216)$$

It follows that

$$I = \frac{1}{2} M a^2. \quad (1.217)$$

Finally, according to Equations (1.212) and (1.217), the angular acceleration of the flywheel is

$$\frac{d\omega}{dt} = \left(\frac{2m}{M+2m} \right) \frac{g}{a}. \quad (1.218)$$

1.7.6 Gravitational Collapse of Star

A star can be thought of as a spherically symmetric body that rotates about an axis passing through its center. The spherical symmetry of the star implies that all three of its principal moments of inertia are equal to one another, and that any axis that passes through the center is a principal axis of rotation. (See Section 1.7.2.)

Suppose that the star in question is rotating about the z -axis. Its moment of inertia is

$$I = \sum_i m_i (x_i^2 + y_i^2), \quad (1.219)$$

where the sum is over all of the mass elements that make up the star. If we model the star as a body of uniform mass density ρ then the previous equation becomes

$$I = \rho \int (x^2 + y^2) dV, \quad (1.220)$$

where dV is a volume element, and the integral is over the volume of the star. In spherical polar coordinates, $x^2 + y^2 = r^2 \sin^2 \theta$ and $dV = r^2 \sin \theta dr d\theta d\phi$. (See Section A.23.) Hence,

$$I = \rho \int_0^a \int_0^\pi \int_0^{2\pi} r^4 \sin^3 \theta dr d\theta d\phi, \quad (1.221)$$

where a is the radius of the star. It follows that

$$I = 2\pi\rho \int_0^a r^4 dr \int_0^\pi \sin^3 \theta d\theta = 2\pi\rho \frac{a^5}{5} \frac{4}{3} = \frac{8\pi\rho a^5}{15}. \quad (1.222)$$

However,

$$M = \frac{4\pi\rho a^3}{3}, \quad (1.223)$$

where M is the mass of the star. The previous two equations yield

$$I = \frac{2}{5} M a^2. \quad (1.224)$$

The angular momentum of the star is

$$L = I \omega, \quad (1.225)$$

where ω is its angular velocity. If the star is isolated, such that it is not subject to an external torque, then its angular equation of motion, (1.180), reduces to

$$\frac{dL}{dt} = 0. \quad (1.226)$$

In other words, the angular momentum of the star is a conserved quantity.

Consider what would happen if a star such as the Sun collapsed gravitationally until it became a neutron star. The radius of the Sun-like star is about 10^6 km. On the other hand, the radius of a neutron star is only about 10 km. Moreover, the Sun rotates at about 1 revolution per month, which corresponds to an angular velocity of 2.5×10^{-6} rad s⁻¹. Thus, we are considering a process in which the star's initial radius and angular velocity are $a_1 = 10^6$ km and $\omega_1 = 2.5 \times 10^{-6}$ rad s⁻¹, respectively, and its final radius is $a_2 = 10$ km. The question is what is the star's final angular velocity, ω_2 . Well, it is clear from Equations (1.224) and (1.225) that if angular momentum is to be conserved during the collapse then we require

$$a_1^2 \omega_1 = a_2^2 \omega_2, \quad (1.227)$$

or

$$\omega_2 = \left(\frac{a_1}{a_2}\right)^2 \omega_1 = \left(\frac{10^6}{10}\right)^2 2.5 \times 10^{-6} = 2.5 \times 10^4 \text{ rad s}^{-1}, \quad (1.228)$$

which corresponds to about 4000 revolutions per second. Thus, we deduce that neutron stars are likely to rotate thousands of times a second (as is indeed the case) as a consequence of the conservation of angular momentum during their formation. Finally, the star's rotational kinetic energy is $K = (1/2)I\omega^2 \propto a^2\omega^2$, so the ratio of the final to the initial kinetic energy of the star is

$$\frac{K_2}{K_1} = \left(\frac{a_2}{a_1}\right)^2 \left(\frac{\omega_2}{\omega_1}\right)^2 = \left(\frac{a_1}{a_2}\right)^2 = 10^{10}. \quad (1.229)$$

1.7.7 Gyroscopic Precession

Consider a gyroscope that consists of a symmetric top of mass M that rotates about its symmetry axis at the constant angular velocity ω . (The rotation axis is obviously a principal axis of rotation.) The top is free to rotate (without friction) about a fixed pivot point P . Suppose that the axis of the top subtends a constant angle θ with the vertical. Let us set up a Cartesian coordinate system such that the z -axis is vertical, and the x - y plane is horizontal. Suppose that the projection of the axis of

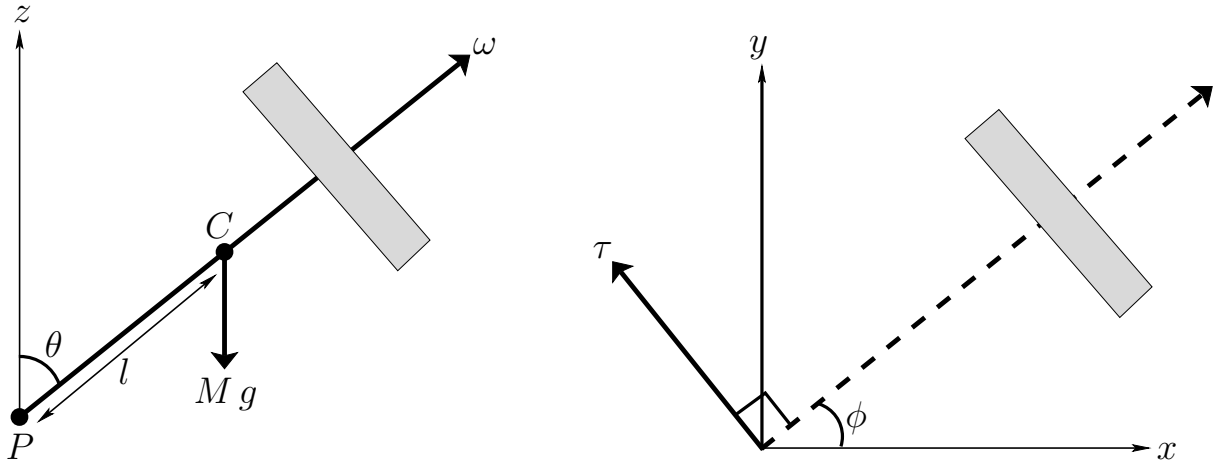


Figure 1.10: A gyroscope.

the top on to the x - y plane subtends an (instantaneous) angle ϕ with the x -axis. Let C be the center of mass of the top (which lies on the symmetry axis), and let the distance PC be l . Finally, let I be the (principal) moment of inertia of the top about its symmetry axis. See Figure 1.10.

It follows, from Figure 1.10, that the angular velocity of the top is

$$\boldsymbol{\omega} = \omega \sin \theta \cos \phi \mathbf{e}_x + \omega \sin \theta \sin \phi \mathbf{e}_y + \omega \cos \theta \mathbf{e}_z. \quad (1.230)$$

Hence, the angular momentum of the top is

$$\mathbf{L} = I \omega \sin \theta \cos \phi \mathbf{e}_x + I \omega \sin \theta \sin \phi \mathbf{e}_y + I \omega \cos \theta \mathbf{e}_z, \quad (1.231)$$

where use has been made of Equation (1.191). Now, the weight of the top exerts a torque $\boldsymbol{\tau}$ that is of magnitude $Mgl \sin \theta$ and whose direction is specified in the figure. Here, g is the acceleration due to gravity. Hence,

$$\boldsymbol{\tau} = -Mgl \sin \theta \sin \phi \mathbf{e}_x + Mgl \sin \theta \cos \phi \mathbf{e}_y. \quad (1.232)$$

The angular equation of motion of the top, (1.180), is

$$\frac{d\mathbf{L}}{dt} = \boldsymbol{\tau}. \quad (1.233)$$

Assuming that θ and ω are constants, whereas ϕ is time-varying, the x -, y -, and z -components of the previous equation are

$$-I \omega \sin \theta \sin \phi \frac{d\phi}{dt} = -Mgl \sin \theta \sin \phi, \quad (1.234)$$

$$I \omega \sin \theta \cos \phi \frac{d\phi}{dt} = Mgl \sin \theta \cos \phi, \quad (1.235)$$

$$0 = 0, \quad (1.236)$$

respectively. Hence, we deduce that the gravitational torque acting on the top causes its axis of rotation to precess about the vertical at the rate

$$\frac{d\phi}{dt} = \frac{M g l}{I \omega}, \quad (1.237)$$

while maintaining a constant inclination to the vertical, and a constant spin rate. Note that the precession is in the same direction as the vertical component of the top's angular velocity. Interestingly enough, the precession rate is independent of the angle of inclination of the rotation axis to the vertical.

1.8 Newtonian Gravity

1.8.1 Gravity

The force that causes objects to fall toward the surface of the Earth, maintains the Moon in orbit about the Earth, and maintains the planets in orbit around the Sun, is called *gravity*, and was first correctly described by Isaac Newton in 1687. According to Newton, any two point mass objects (or spherically symmetric objects of finite extent) exert a force of attraction on one another. This force points along the line of centers joining the objects, is directly proportional to the product of the objects' masses, and inversely proportional to the square of the distance between them.

Consider a system consisting of two point mass objects. Let object 1 have mass m_1 and displacement \mathbf{r}_1 . Let object 2 have mass m_2 and displacement \mathbf{r}_2 . The gravitational force exerted on object 2 by object 1 is written

$$\mathbf{f}_{21} = -G m_1 m_2 \frac{\mathbf{r}_2 - \mathbf{r}_1}{|\mathbf{r}_2 - \mathbf{r}_1|^3}. \quad (1.238)$$

The constant of proportionality, G , is called the *universal gravitational constant*, and takes the value

$$G = 6.67430 \times 10^{-11} \text{ m}^3 \text{ kg}^{-1} \text{ s}^{-2}. \quad (1.239)$$

This constant was first 'measured' by Henry Cavendish in 1798 (to be more exact, the result $G = 6.74 \times 10^{-11} \text{ m}^3 \text{ kg}^{-1} \text{ s}^{-2}$ can be inferred from Cavendish's results). An equal and opposite force to (1.238) acts on object 1.

Suppose that we have a system of N point mass objects. Let the i th object have mass m_i and displacement \mathbf{r}_i . Now, it is an experimentally verified fact that gravity is a *superposable* force. In other words, the gravitational force exerted on object i by object j is unaffected by the presence of any other objects in the universe. Hence, the net gravitational force experienced by object i is

$$\mathbf{f}_i = -G m_i \sum_{j=1, N}^{j \neq i} m_j \frac{\mathbf{r}_i - \mathbf{r}_j}{|\mathbf{r}_i - \mathbf{r}_j|^3}. \quad (1.240)$$

Note that object i is missing from the sum on the right-hand side of the previous equation because this object cannot exert a gravitational force on itself.

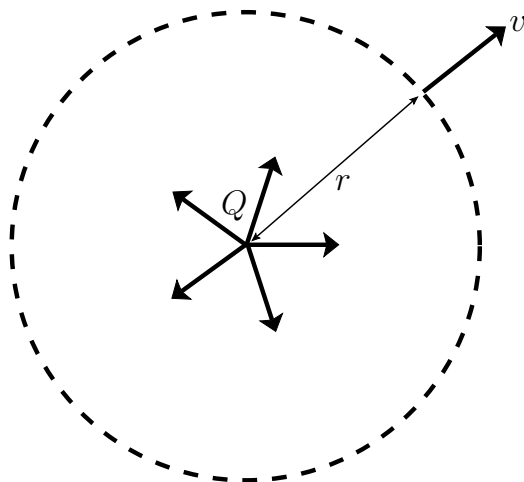


Figure 1.11: Fluid outflow from a point source.

Suppose, finally, that a point object of mass m is located at the origin of our coordinate system. It follows from Equation (1.238) that the gravitational acceleration, due to the gravitational attraction of mass m , experienced by another point object whose displacement is \mathbf{r} is

$$\mathbf{g} = -G m \frac{\mathbf{r}}{r^3}. \quad (1.241)$$

1.8.2 Gauss's Law

Gauss's law applies to any inverse-square central force. This law was first discovered by Joseph-Louis Lagrange in 1773, and was later rediscovered by Carl Friedrich Gauss in 1835.

Suppose that we have an isolated point source that emits an incompressible fluid isotropically in all directions at the rate Q ($\text{m}^3 \text{s}^{-1}$). By symmetry, we would expect the fluid to flow radially away from the source, isotropically in all directions. In other words, if the source is located at the origin then we expect the fluid velocity at displacement \mathbf{r} to be of the form

$$\mathbf{v} = v(r) \frac{\mathbf{r}}{r}. \quad (1.242)$$

Let us surround our source by an imaginary spherical surface of radius r . See Figure 1.11. The net volume rate of flow of fluid out of the surface is $4\pi r^2 v(r)$. However, if the fluid is incompressible (and the flow pattern has achieved a steady-state) then the volume rate of flow out of the surface must equal the volume rate of flow from the source (otherwise, the fluid inside the surface would suffer compression or rarefaction). In other words, $4\pi r^2 v(r) = Q$. Hence, Equation (1.242) becomes

$$\mathbf{v} = \frac{Q}{4\pi} \frac{\mathbf{r}}{r^3}. \quad (1.243)$$

It can be seen that the previous expression has an analogous form to expression (1.241), provided that we make the identifications $\mathbf{v} \rightarrow \mathbf{g}$ and $Q \rightarrow -4\pi G m$.

Suppose that we have N point sources of incompressible fluid. Let the i th source have a volume rate of flow Q_i . Let us surround these sources by an imaginary closed surface (that is not necessarily spherical), S . Now, the volume rate of flow of fluid out of S is $\oint_S \mathbf{v} \cdot d\mathbf{S}$, where $\mathbf{v}(\mathbf{r})$ is the flow field, and $d\mathbf{S}$ is an (outward pointing) element of S . (See Section A.16.) However, if the fluid is incompressible (and the flow pattern has achieved a steady-state) then the volume rate of flow out of S must match the sum of the volume rates of flow of the sources within S (otherwise, the fluid inside the surface would suffer compression or rarefaction). Thus, we deduce that

$$\oint_S \mathbf{v} \cdot d\mathbf{S} = \sum_{i=1,N} Q_i. \quad (1.244)$$

Now, if there are sources that lie outside S then they do not affect the previous relation (because a source outside S gives rise to zero net flow of incompressible fluid out of S). Hence, we can interpret the sum in the previous equation as a sum that includes all sources that lie inside S , but excludes any sources that lie outside S .

Finally, we can exploit the previously mentioned analogy between incompressible fluid flow and gravitational acceleration to deduce the following result. Suppose that there are N point objects of mass m_i . Let us surround these objects by an imaginary surface S . Making use of the identifications $\mathbf{v} \rightarrow \mathbf{g}$ and $Q_i \rightarrow -4\pi G m_i$, the previous equation transforms to give

$$\oint_S \mathbf{g} \cdot d\mathbf{S} = -4\pi G \sum_{i=1,N} m_i. \quad (1.245)$$

As before, if there are objects outside S then they do not affect the previous relation. Thus, we deduce that the flux of gravitational acceleration out of an arbitrary closed surface, S , is equal to $-4\pi G$ multiplied by the sum of the masses of any objects lying inside the surface. This is Gauss's law. The imaginary surface S is known as a *Gaussian surface*.

Suppose that, instead of having a collection of point objects, we have a continuous mass distribution whose mass density is $\rho(\mathbf{r})$. The previous equation generalizes to give

$$\oint_S \mathbf{g} \cdot d\mathbf{S} = -4\pi G \int_V \rho dV, \quad (1.246)$$

where V is the volume enclosed by S , and dV is an element of V . (See Section A.17.) Now, according to the divergence theorem (see Section A.20),

$$\oint_S \mathbf{g} \cdot d\mathbf{S} = \int_V \nabla \cdot \mathbf{g} dV, \quad (1.247)$$

where $\nabla \cdot \mathbf{g} = \partial g_x / \partial x + \partial g_y / \partial y + \partial g_z / \partial z$ is the *divergence* of the acceleration field. The previous two equations yield

$$\int_V \nabla \cdot \mathbf{g} dV = -4\pi G \int_V \rho dV. \quad (1.248)$$

However, the volume V is arbitrary, so the only way that the previous equation could hold for all possible volumes is if

$$\nabla \cdot \mathbf{g} = -4\pi G \rho. \quad (1.249)$$

This is the differential form of Gauss's law.

1.8.3 Gravitational Field of Earth

Suppose that we model the Earth (not very accurately) as a uniform sphere of mass M and radius R . Let the center of the Earth lie at the origin. We wish to determine the gravitational acceleration due to the Earth, \mathbf{g} , both inside and outside the Earth. By symmetry, we expect the gravitational acceleration at a general point whose displacement is \mathbf{r} to be directed radially inward toward the center of the Earth, and to only depend on the radial distance, r , from the center of the Earth. In other words,

$$\mathbf{g} = -a(r) \frac{\mathbf{r}}{r}. \quad (1.250)$$

Consider a spherical Gaussian surface, S , of radius r , whose center corresponds to the center of the Earth. The flux of gravitational acceleration out of this surface is

$$\oint_S \mathbf{g} \cdot d\mathbf{S} = -4\pi r^2 a(r). \quad (1.251)$$

Thus, Gauss's law, (1.246), yields

$$a(r) = \frac{G m(r)}{r^2}, \quad (1.252)$$

where $m(r)$ is the mass enclosed by the surface. Now, if $r < R$ then $m(r) = (r/R)^3 M$ (by proportion), but if $r > R$ then $m(r) = M$. Thus, we deduce that

$$a(r) = \frac{G M r}{R^3} \quad (1.253)$$

for $r < R$, and

$$a(r) = \frac{G M}{r^2} \quad (1.254)$$

for $r > R$. In other words, inside the Earth, the gravitational acceleration due to the Earth increases linearly with distance from the Earth's center, but, outside the Earth, it falls off as the inverse-square of distance from the Earth's center. In particular, the gravitational field outside the Earth is exactly the same as that of a point object whose mass is equal to that of the Earth, and that is located at the Earth's center. This is an example of an important result first derived by Newton; namely, that the gravitational field outside a spherically symmetric mass distribution is the same as that of a point object located at the center of the distribution whose mass is equal to that of the distribution.

It is clear from Equation (1.254) that the gravitational acceleration at the surface of the Earth is

$$g \equiv a(R) = \frac{G M}{R^2}. \quad (1.255)$$

Given that the measured (average) gravitational acceleration at the Earth's surface is $g = 9.81 \text{ m s}^{-2}$, and that the measured (mean) radius of the Earth is $R = 6.371 \times 10^6 \text{ m}$, we arrive at the following estimate for the Earth's mass,

$$M = \frac{g R^2}{G} = \frac{9.81 \times (6.371 \times 10^6)^2}{6.67430 \times 10^{-11}} = 5.965 \times 10^{24} \text{ kg} \quad (1.256)$$

This estimate lies within 0.1% of the correct value, which is $M = 5.972 \times 10^{24} \text{ kg}$

1.8.4 Gravitational Potential Energy

Suppose that a spherically symmetric object of mass M is located at the origin of our coordinate system. The gravitational force, due to the gravitational attraction of mass M , experienced by a point object of mass m and displacement is \mathbf{r} (located outside the former mass) is

$$\mathbf{f} = -G M m \frac{\mathbf{r}}{r^3}. \quad (1.257)$$

[See Equation (1.241).] Now, $r = (x^2 + y^2 + z^2)^{1/2}$. It is easily demonstrated that

$$\frac{\partial r}{\partial x} = \frac{x}{r}, \quad (1.258)$$

$$\frac{\partial r}{\partial y} = \frac{y}{r}, \quad (1.259)$$

$$\frac{\partial r}{\partial z} = \frac{z}{r}. \quad (1.260)$$

Consider

$$\nabla \left(\frac{1}{r} \right) \equiv \frac{\partial(1/r)}{\partial x} \mathbf{e}_x + \frac{\partial(1/r)}{\partial y} \mathbf{e}_y + \frac{\partial(1/r)}{\partial z} \mathbf{e}_z. \quad (1.261)$$

(See Section A.19.) It follows that

$$\begin{aligned} \nabla \left(\frac{1}{r} \right) &= -\frac{1}{r^2} \frac{\partial r}{\partial x} \mathbf{e}_x - \frac{1}{r^2} \frac{\partial r}{\partial y} \mathbf{e}_y - \frac{1}{r^2} \frac{\partial r}{\partial z} \mathbf{e}_z \\ &= -\frac{x}{r^3} \mathbf{e}_x - \frac{y}{r^3} \mathbf{e}_y - \frac{z}{r^3} \mathbf{e}_z \\ &= -\frac{\mathbf{r}}{r^3}, \end{aligned} \quad (1.262)$$

where use has been made of Equations (1.258)–(1.260). The previous equation can be combined with Equation (1.257) to give

$$\mathbf{f} = \nabla \left(\frac{G M m}{r} \right). \quad (1.263)$$

A comparison with Equation (1.47) reveals that the gravitational force field of our spherical object is a conservative field with the associated potential energy

$$U(\mathbf{r}) = -\frac{G M m}{r}. \quad (1.264)$$

Note that, by convention, the potential energy at infinity is zero.

A particle of mass m moving in the gravitational field of our spherical object has a conserved energy

$$E = K + U = \frac{1}{2} m v^2 - \frac{G M m}{r}, \quad (1.265)$$

where v is the particle's instantaneous speed. (See Sections 1.3.2 and 1.3.5.)

Let us again model the Earth as a sphere of mass M and radius R that is centered at the origin. Consider an object that is launched from the surface of the Earth, in an arbitrary outward direction, with speed v_{escape} . Suppose that the object only just manages to escape from the Earth's gravitational field. It follows that the object's speed at infinity (i.e., $1/r = 0$) is zero. Thus, it is clear from the previous equation that the object's conserved energy, E , is also zero. Hence,

$$0 = \frac{1}{2} m v_{\text{escape}}^2 - \frac{G M m}{R} \quad (1.266)$$

at the surface of the Earth, which implies that

$$v_{\text{escape}} = \left(\frac{2 G M}{R} \right)^{1/2} = \left[\frac{2 \times (6.67430 \times 10^{-11}) \times (5.972 \times 10^{24})}{6.371 \times 10^6} \right]^{1/2} = 11.19 \text{ km s}^{-1}. \quad (1.267)$$

The speed v_{escape} , which is known as the *escape speed*, is the minimum speed at which an object must be launched from the Earth's surface if it is to reach outer space.

1.8.5 Gravitational Potential

We have seen that the force experienced by a point object of mass m situated in a gravitational field can be written

$$\mathbf{f} = -\nabla U, \quad (1.268)$$

where U is the object's gravitational potential energy. It is clear from Equation (1.264) that $U \propto m$. It follows that the gravitational acceleration of the object, $\mathbf{g} = \mathbf{f}/m$, can be written

$$\mathbf{g} = -\nabla \Phi, \quad (1.269)$$

where $\Phi = U/m$ is independent of m . The quantity Φ is known as *gravitational potential*.

From Equation (1.264), the gravitational potential due to a point object (or a spherically symmetric object) of mass M situated at the origin is

$$\Phi(\mathbf{r}) = -\frac{G M}{r}. \quad (1.270)$$

Consider a collection of N point objects. Let the i th object have mass m_i and displacement \mathbf{r}_i . Given that gravity is a superposable force, the generalization of the previous equation is clearly

$$\Phi(\mathbf{r}) = -\sum_{i=1, N} \frac{G m_i}{|\mathbf{r} - \mathbf{r}_i|}. \quad (1.271)$$

Moreover, the gravitational acceleration due to the collection of objects is again given by Equation (1.269). Finally, if, instead of having a collection of point mass objects, we have a continuous mass distribution characterized by a mass density $\rho(\mathbf{r})$, then the previous expression generalizes to give

$$\Phi(\mathbf{r}) = -\int \frac{G \rho(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} dV', \quad (1.272)$$

where the integral is over all space.

Equation (1.269) can be combined with the differential form of Gauss's theorem, (1.249), to give

$$\nabla^2 \Phi = 4\pi G \rho. \quad (1.273)$$

Here,

$$\nabla^2 \Phi \equiv \nabla \cdot \nabla \Phi = \frac{\partial^2 \Phi}{\partial x^2} + \frac{\partial^2 \Phi}{\partial y^2} + \frac{\partial^2 \Phi}{\partial z^2}, \quad (1.274)$$

is known as the *Laplacian* of Φ . (See Section A.21.) Equation (1.273), which specifies the gravitational potential, $\Phi(\mathbf{r})$, generated by a continuous mass distribution of mass density $\rho(\mathbf{r})$, is known as *Poisson's equation*. Of course, Equation (1.272) is the integral form of Poisson's equation.

1.9 Planetary Motion

1.9.1 Kepler's Laws

As is well known, Johannes Kepler was the first astronomer to correctly describe planetary motion in the solar system (in works published between 1609 and 1619). The motion of the planets is summed up in three simple laws:

1. The planetary orbits are all ellipses that are confocal with the Sun (i.e., the Sun lies at one of the foci of each ellipse).
2. The radius vectors connecting each planet to the Sun sweep out equal areas in equal time intervals.
3. The squares of the orbital periods of the planets are proportional to the cubes of their orbital major radii.

Let us now see if we can derive Kepler's laws from Newton's laws of motion.

Suppose that the Sun, which is of mass M , is located at the origin of our coordinate system. Consider a planet, of mass m , whose instantaneous displacement is \mathbf{r} . The gravitational force exerted on the planet by the Sun is thus written

$$\mathbf{f} = -\frac{GMm}{r^3} \mathbf{r}. \quad (1.275)$$

[See Equation (1.257).] An equal and opposite force to (1.275) acts on the Sun. However, we shall assume that the Sun is so much more massive than the planet in question that this force does not cause the Sun's position to shift appreciably. Hence, the Sun will always remain at the origin of our coordinate system. Likewise, we shall neglect the gravitational forces exerted on our planet by the other planets in the solar system, compared to the much larger gravitational force exerted by the Sun. Thus, according to Newton's second law, the equation of motion of our planet is

$$m \frac{d^2 \mathbf{r}}{dt^2} = \mathbf{f}, \quad (1.276)$$

which reduces to

$$\frac{d^2 \mathbf{r}}{dt^2} = -\frac{GM}{r^3} \mathbf{r}. \quad (1.277)$$

Note that the planetary mass, m , has cancelled out on both sides of the previous equation.

1.9.2 Planetary Conservation Laws

As we have seen, gravity is a conservative force. Hence, the gravitational force (1.275) can be written

$$\mathbf{f} = -\nabla U, \quad (1.278)$$

where the potential energy, $U(\mathbf{r})$, of our planet in the Sun's gravitational field takes the form

$$U(\mathbf{r}) = -\frac{GMm}{r}. \quad (1.279)$$

[See Equation (1.264).] It follows that the total energy of our planet is a conserved quantity. (See Section 1.3.5.) In other words,

$$\mathcal{E} = \frac{v^2}{2} - \frac{GM}{r} \quad (1.280)$$

is constant in time. Here, \mathcal{E} is actually the planet's total energy per unit mass, and $\mathbf{v} = d\mathbf{r}/dt$.

Gravity is also a central force. This means that the gravitational force exerted on our planet is always directed toward the origin of our coordinate system (i.e., the Sun), which implies that the force exerts zero torque about the origin. Hence, the angular momentum of our planet (about the origin) is a conserved quantity. (See Section 1.4.5.) In other words,

$$\mathbf{h} = \mathbf{r} \times \mathbf{v}, \quad (1.281)$$

which is actually the planet's angular momentum per unit mass, is constant in time. Taking the scalar product of the previous equation with \mathbf{r} , we obtain

$$\mathbf{h} \cdot \mathbf{r} = 0. \quad (1.282)$$

This is the equation of a plane that passes through the origin, and whose normal is parallel to \mathbf{h} . Because \mathbf{h} is a constant vector, it always points in the same direction. We, therefore, conclude that the motion of our planet is two-dimensional in nature; that is, it is confined to some fixed plane that passes through the origin. Without loss of generality, we can let this plane coincide with the x - y plane.

1.9.3 Plane Polar Coordinates

We can determine the instantaneous position of our planet in the x - y plane in terms of standard Cartesian coordinates, (x, y) , or plane polar coordinates, (r, θ) , as illustrated in Figure 1.12. Here, $r = (x^2 + y^2)^{1/2}$ and $\theta = \tan^{-1}(y/x)$. It is helpful to define two unit vectors, $\mathbf{e}_r \equiv \mathbf{r}/r$ and $\mathbf{e}_\theta \equiv \mathbf{e}_z \times \mathbf{e}_r$, at the instantaneous position of the planet. The first always points radially away from the origin,

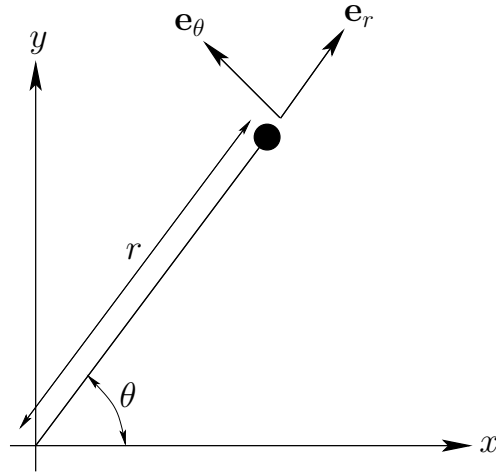


Figure 1.12: Plane polar coordinates.

whereas the second is normal to the first, in the direction of increasing θ . As is easily demonstrated, the Cartesian components of \mathbf{e}_r and \mathbf{e}_θ are

$$\mathbf{e}_r = (\cos \theta, \sin \theta), \quad (1.283)$$

$$\mathbf{e}_\theta = (-\sin \theta, \cos \theta), \quad (1.284)$$

respectively.

We can write the displacement of our planet as

$$\mathbf{r} = r \mathbf{e}_r. \quad (1.285)$$

Thus, the planet's velocity becomes

$$\mathbf{v} = \frac{d\mathbf{r}}{dt} = \dot{r} \mathbf{e}_r + r \dot{\mathbf{e}}_r, \quad (1.286)$$

where $\dot{}$ is shorthand for d/dt . Note that \mathbf{e}_r has a non-zero time-derivative (unlike a Cartesian unit vector) because its direction changes as the planet moves around. As is easily demonstrated, from differentiating Equation (1.283) with respect to time,

$$\dot{\mathbf{e}}_r = \dot{\theta}(-\sin \theta, \cos \theta) = \dot{\theta} \mathbf{e}_\theta. \quad (1.287)$$

Thus,

$$\mathbf{v} = \dot{r} \mathbf{e}_r + r \dot{\theta} \mathbf{e}_\theta. \quad (1.288)$$

The planet's acceleration is written

$$\mathbf{a} = \frac{d\mathbf{v}}{dt} = \frac{d^2\mathbf{r}}{dt^2} = \ddot{r} \mathbf{e}_r + \dot{r} \dot{\mathbf{e}}_r + (\dot{r} \dot{\theta} + r \ddot{\theta}) \mathbf{e}_\theta + r \dot{\theta} \dot{\mathbf{e}}_\theta. \quad (1.289)$$

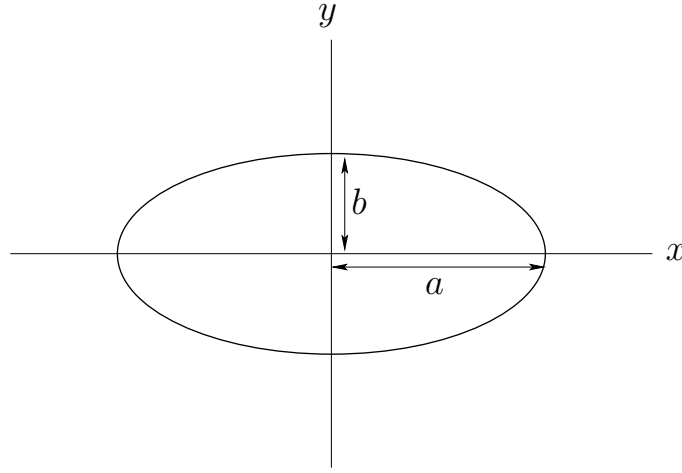


Figure 1.13: An ellipse.

Again, \mathbf{e}_θ has a non-zero time-derivative because its direction changes as the planet moves around. Differentiation of Equation (1.284) with respect to time yields

$$\dot{\mathbf{e}}_\theta = \dot{\theta}(-\cos \theta, -\sin \theta) = -\dot{\theta} \mathbf{e}_r. \quad (1.290)$$

Hence,

$$\mathbf{a} = (\ddot{r} - r\dot{\theta}^2) \mathbf{e}_r + (r\ddot{\theta} + 2\dot{r}\dot{\theta}) \mathbf{e}_\theta. \quad (1.291)$$

It follows that the equation of motion of our planet, (1.277), can be written

$$\mathbf{a} = (\ddot{r} - r\dot{\theta}^2) \mathbf{e}_r + (r\ddot{\theta} + 2\dot{r}\dot{\theta}) \mathbf{e}_\theta = -\frac{GM}{r^2} \mathbf{e}_r. \quad (1.292)$$

Because \mathbf{e}_r and \mathbf{e}_θ are mutually orthogonal, we can separately equate the coefficients of both, in the previous equation, to give a radial equation of motion,

$$\ddot{r} - r\dot{\theta}^2 = -\frac{GM}{r^2}, \quad (1.293)$$

and a tangential equation of motion,

$$r\ddot{\theta} + 2\dot{r}\dot{\theta} = 0. \quad (1.294)$$

1.9.4 Conic Sections

The ellipse, the parabola, and the hyperbola are collectively known as *conic sections*, because these three types of curve can be obtained by taking various different plane sections of a right cone. It turns out that the possible solutions of Equations (1.293) and (1.294) are all conic sections. It is, therefore, appropriate for us to briefly review these curves.

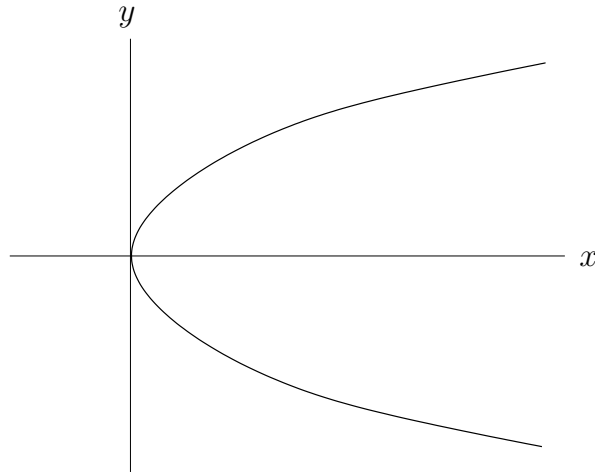


Figure 1.14: A parabola.

An *ellipse*, centered on the origin, of major radius a and minor radius b , which are aligned along the x - and y -axes, respectively (see Figure 1.13), satisfies the following well-known equation:

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1. \quad (1.295)$$

Likewise, a *parabola* that is aligned along the $+x$ -axis, and passes through the origin (see Figure 1.14), satisfies:

$$y^2 - b x = 0, \quad (1.296)$$

where $b > 0$.

Finally, a *hyperbola* that is aligned along the $+x$ -axis, and whose asymptotes intersect at the origin (see Figure 1.15), satisfies:

$$\frac{x^2}{a^2} - \frac{y^2}{b^2} = 1. \quad (1.297)$$

Here, a is the distance of closest approach to the origin. The asymptotes subtend an angle $\phi = \tan^{-1}(b/a)$ with the x -axis.

It is not clear, at this stage, what the ellipse, the parabola, and the hyperbola have in common (other than being conic sections). Well, it turns out that what these three curves have in common is that they can all be represented as the locus of a movable point whose distance from a fixed point is in a constant ratio to its perpendicular distance to some fixed straight-line. Let the fixed point (which is termed the *focus* of the ellipse/parabola/hyperbola) lie at the origin, and let the fixed line correspond to $x = -d$ (with $d > 0$). Thus, the distance of a general point (x, y) (which lies to the right of the line $x = -d$) from the origin is $r_1 = (x^2 + y^2)^{1/2}$, whereas the perpendicular distance of the point from the line $x = -d$ is $r_2 = x + d$. See Figure 1.16. In polar coordinates, $r_1 = r$ and $r_2 = r \cos \theta + d$. Hence, the locus of a point for which r_1 and r_2 are in a fixed ratio satisfies the following equation:

$$\frac{r_1}{r_2} = \frac{\sqrt{x^2 + y^2}}{x + d} = \frac{r}{r \cos \theta + d} = e, \quad (1.298)$$

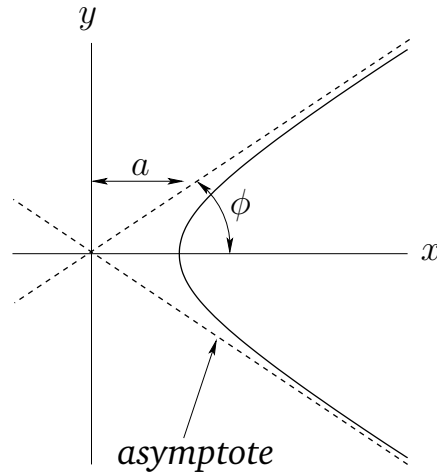


Figure 1.15: A hyperbola.

where $e \geq 0$ is a constant. When expressed in terms of polar coordinates, the previous equation can be rearranged to give

$$r = \frac{r_c}{1 - e \cos \theta}, \quad (1.299)$$

where $r_c = e d$.

When written in terms of Cartesian coordinates, Equation (1.298) can be rearranged to give

$$\frac{(x - x_c)^2}{a^2} + \frac{y^2}{b^2} = 1, \quad (1.300)$$

for $e < 1$. Here,

$$a = \frac{r_c}{1 - e^2}, \quad (1.301)$$

$$b = \frac{r_c}{\sqrt{1 - e^2}} = \sqrt{1 - e^2} a, \quad (1.302)$$

$$x_c = \frac{e r_c}{1 - e^2} = e a. \quad (1.303)$$

Equation (1.300) can be recognized as the equation of an ellipse whose center lies at $(x_c, 0)$, and whose major and minor radii, a and b , are aligned along the x - and y -axes, respectively [cf., Equation (1.295)].

When again written in terms of Cartesian coordinates, Equation (1.298) can be rearranged to give

$$y^2 - 2 r_c (x - x_c) = 0, \quad (1.304)$$

for $e = 1$. Here, $x_c = -r_c/2$. This is the equation of a parabola that passes through the point $(x_c, 0)$, and that is aligned along the $+x$ -direction [cf., Equation (1.296)].

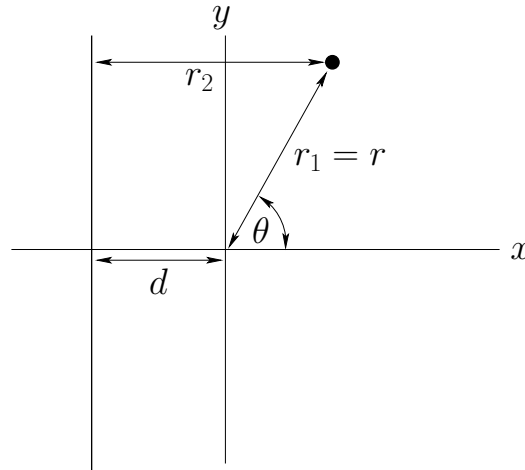


Figure 1.16: Conic sections in plane polar coordinates.

Finally, when written in terms of Cartesian coordinates, Equation (1.298) can be rearranged to give

$$\frac{(x - x_c)^2}{a^2} - \frac{y^2}{b^2} = 1, \quad (1.305)$$

for $e > 1$. Here,

$$a = \frac{r_c}{e^2 - 1}, \quad (1.306)$$

$$b = \frac{r_c}{\sqrt{e^2 - 1}} = \sqrt{e^2 - 1} a, \quad (1.307)$$

$$x_c = -\frac{e r_c}{e^2 - 1} = -e a. \quad (1.308)$$

Equation (1.305) can be recognized as the equation of a hyperbola whose asymptotes intersect at $(x_c, 0)$, and that is aligned along the $+x$ -direction [cf., Equation (1.297)]. The asymptotes subtend an angle

$$\phi = \tan^{-1} \left(\frac{b}{a} \right) = \tan^{-1} \left(\sqrt{e^2 - 1} \right) \quad (1.309)$$

with the x -axis.

In conclusion, Equation (1.299) is the polar equation of a general conic section that is confocal with the origin. For $e < 1$, the conic section is an ellipse. For $e = 1$, the conic section is a parabola. Finally, for $e > 1$, the conic section is a hyperbola.

1.9.5 Kepler's Second Law

Multiplying our planet's tangential equation of motion, (1.294), by r , we obtain

$$r^2 \ddot{\theta} + 2r \dot{r} \dot{\theta} = 0. \quad (1.310)$$

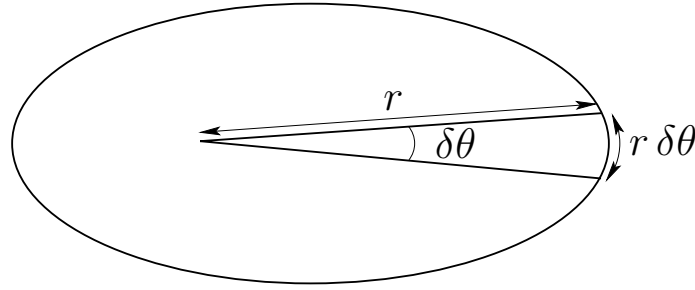


Figure 1.17: Kepler's second law.

However, the previous equation can be also written

$$\frac{d(r^2 \dot{\theta})}{dt} = 0, \quad (1.311)$$

which implies that

$$h = r^2 \dot{\theta} \quad (1.312)$$

is constant in time. It is easily demonstrated that h is the magnitude of the vector \mathbf{h} defined in Equation (1.281). Thus, the fact that h is constant in time is equivalent to the statement that the angular momentum of our planet is a constant of its motion. As we have already mentioned, this is the case because gravity is a central force.

Suppose that the radius vector connecting our planet to the origin (i.e., the Sun) sweeps out an angle $\delta\theta$ between times t and $t + \delta t$. See Figure 1.17. The approximately triangular region swept out by the radius vector has the area

$$\delta A \simeq \frac{1}{2} r^2 \delta\theta, \quad (1.313)$$

because the area of a triangle is half its base ($r \delta\theta$) times its height (r). Hence, the rate at which the radius vector sweeps out area is

$$\frac{dA}{dt} = \lim_{\delta t \rightarrow 0} \frac{r^2 \delta\theta}{2 \delta t} = \frac{r^2}{2} \frac{d\theta}{dt} = \frac{h}{2}. \quad (1.314)$$

Thus, the radius vector sweeps out area at a constant rate (because h is constant in time). This is Kepler's second law. We conclude that Kepler's second law of planetary motion is a direct consequence of angular momentum conservation.

1.9.6 Kepler's First Law

Our planet's radial equation of motion, (1.293), can be combined with Equation (1.312) to give

$$\ddot{r} - \frac{h^2}{r^3} = -\frac{GM}{r^2}. \quad (1.315)$$

Suppose that $r = u^{-1}$, where $u = u(\theta)$. It follows that

$$\dot{r} = -\frac{\dot{u}}{u^2} = -r^2 \frac{du}{d\theta} \frac{d\theta}{dt} = -h \frac{du}{d\theta}. \quad (1.316)$$

Likewise,

$$\ddot{r} = -h \frac{d^2u}{d\theta^2} \dot{\theta} = -u^2 h^2 \frac{d^2u}{d\theta^2}. \quad (1.317)$$

Hence, Equation (1.315) can be written in the linear form

$$\frac{d^2u}{d\theta^2} + u = \frac{GM}{h^2}. \quad (1.318)$$

The general solution to the previous equation is

$$u(\theta) = \frac{GM}{h^2} [1 - e \cos(\theta - \theta_0)], \quad (1.319)$$

where e and θ_0 are arbitrary constants. Without loss of generality, we can set $\theta_0 = 0$ by rotating our coordinate system about the z -axis. Thus, we obtain

$$r(\theta) = \frac{r_c}{1 - e \cos \theta}, \quad (1.320)$$

where

$$r_c = \frac{h^2}{GM}. \quad (1.321)$$

We immediately recognize Equation (1.320) as the equation of a conic section that is confocal with the origin (i.e., with the Sun). Specifically, for $e < 1$, Equation (1.320) is the equation of an ellipse that is confocal with the Sun. Thus, the orbit of our planet around the Sun is a confocal ellipse. This is Kepler's first law of planetary motion. Of course, a planet cannot have a parabolic or a hyperbolic orbit, because such orbits are only appropriate to objects that are ultimately able to escape from the Sun's gravitational field.

For the case of an elliptic orbit, the *eccentricity*, e , measures the displacement of the Sun from the geometric center of the orbit; in fact, according to Equation (1.303), this displacement is $e a$, where $a = r_c/(e^2 - 1)$ is the major radius. The eccentricity also measures the elongation of the orbit; in fact, according to Equations (1.301) and (1.302), $b/a = \sqrt{1 - e^2}$, where b is the minor radius. Note that the displacement is first order in e , whereas the elongation is second order. As is clear from Table 1.4, the planets in the solar system all have orbits characterized by small eccentricities, which implies that these orbits are actually all quite close to being circular.

1.9.7 Kepler's Third Law

We have seen that the radius vector connecting our planet to the origin sweeps out area at the constant rate $dA/dt = h/2$. [See Equation (1.314).] We have also seen that the planetary orbit is an ellipse. Suppose that the major and minor radii of the ellipse are a and b , respectively. It follows

Planet	$a(\text{AU})$	e	$T(\text{yr})$	T^2/a^3
Mercury	0.3871	0.20564	0.241	1.0013
Venus	0.7233	0.00676	0.615	0.9995
Earth	1.0000	0.01673	1.000	1.0000
Mars	1.5237	0.09337	1.881	1.0002
Jupiter	5.2025	0.04854	11.87	1.0006
Saturn	9.5415	0.05551	29.47	0.9998
Uranus	19.188	0.04686	84.05	1.0000
Neptune	30.070	0.00895	164.9	1.0001

Table 1.4: Orbits of planets in the solar system. a - major radius; e - eccentricity; T - orbital period. Here, an astronomical unit (AU) is 1.496×10^{11} m.

that the area of the ellipse is $A = \pi a b$. Now, we expect the radius vector to sweep out the whole area of the ellipse in a single orbital period, T . Hence,

$$T = \frac{A}{(dA/dt)} = \frac{2\pi a b}{h}. \quad (1.322)$$

It follows from Equations (1.301), (1.302), and (1.321) that

$$T^2 = \frac{4\pi^2 a^3}{GM}. \quad (1.323)$$

In other words, the square of the orbital period of our planet is proportional to the cube of its orbital major radius. This is Kepler's third law of planetary motion. As is clear from Table 1.4, Kepler's third law very accurately describes the orbits of the planets in the solar system.

1.9.8 Orbital Energies

According to Equations (1.320) and (1.321), when $\theta = \pi$ an object moving in the Sun's gravitational attains its closest distance to the Sun,

$$r_p = \frac{h^2}{GM(1+e)}. \quad (1.324)$$

This distance is known as the *perihelion* distance. At the point of closest approach to the Sun, the object's instantaneous radial velocity, \dot{r} , is zero (because r attains a minimum value at this point). Hence, making use of Equations (1.288) and (1.312), the object's speed at the perihelion distance is

$$v^2 = r_p^2 \dot{\theta}^2 = \frac{h^2}{r_p^2}. \quad (1.325)$$

Thus, according to Equation (1.280) and (1.324), the object's energy per unit mass at the perihelion distance is

$$\mathcal{E} = \frac{h^2}{2r_p^2} - \frac{G}{r_p} = \frac{GM(1+e)}{2r_p} - \frac{GM}{r_p}, \quad (1.326)$$

which reduces to

$$\mathcal{E} = \frac{GM}{2r_p}(e - 1). \quad (1.327)$$

Of course, because \mathcal{E} is a conserved quantity, the previous expression specifies the energy per unit mass of the object at all distances from the Sun. We conclude that an object in an elliptic orbit ($e < 1$) has a negative total energy, whereas an object in a parabolic orbit ($e = 1$) has zero total energy, and an object in a hyperbolic orbit ($e > 1$) has a positive total energy. This makes sense, because in a conservative system in which the potential energy at infinity is set to zero [see Equation (1.279)], we expect bounded orbits to have negative total energies, and unbounded orbits to have positive total energies. (See Section 1.3.6.) Thus, elliptical orbits, which are clearly bounded, should indeed have negative total energies, whereas hyperbolic orbits, which are clearly unbounded, should indeed have positive total energies. Parabolic orbits are marginally bounded (i.e., an object executing a parabolic orbit only just escapes from the Sun's gravitational field), and thus have zero total energy.

1.10 Spheroidal Mass Distributions

1.10.1 Gravitational Potential of Uniform Spheroid

Let us use Poisson's equation, (1.273), to calculate the gravitational potential generated around a *spheroid* of uniform mass density γ and mean radius R . A spheroid is the solid body produced by rotating an ellipse about a major or a minor axis. Let the center of the spheroid be located at the origin, let its axis of rotation coincide with the z -axis, and let its outer boundary satisfy

$$r = R_\theta(\theta) = R \left[1 - \frac{2}{3} \epsilon P_2(\cos \theta) \right], \quad (1.328)$$

where ϵ is termed the *ellipticity*. Here, r , θ , ϕ are conventional spherical coordinates. (See Section A.23.) Moreover,

$$P_2(x) = \frac{1}{2} (3x^2 - 1) \quad (1.329)$$

is a *Legendre polynomial* of degree 2. It can be seen that the radius of the spheroid at the poles (i.e., along the rotation axis, $\theta = 0$) is $R_p = R(1 - 2\epsilon/3)$, whereas the radius at the equator (i.e., in the bisecting plane perpendicular to the axis, $\theta = \pi/2$) is $R_e = R(1 + \epsilon/3)$. Hence,

$$\epsilon = \frac{R_e - R_p}{R}. \quad (1.330)$$

Let us assume that $|\epsilon| \ll 1$, so that the spheroid is very close to being a sphere. If $\epsilon > 0$ then the spheroid is slightly squashed along its axis of rotation, and is termed *oblate*. Likewise, if $\epsilon < 0$ then the spheroid is slightly elongated along its axis, and is termed *prolate*. See Figure 1.18. Of course, if $\epsilon = 0$ then the spheroid reduces to a sphere. Note that R is the surface-averaged radius of the spheroid, which implies that the volume of the spheroid is equal to that of a sphere of radius

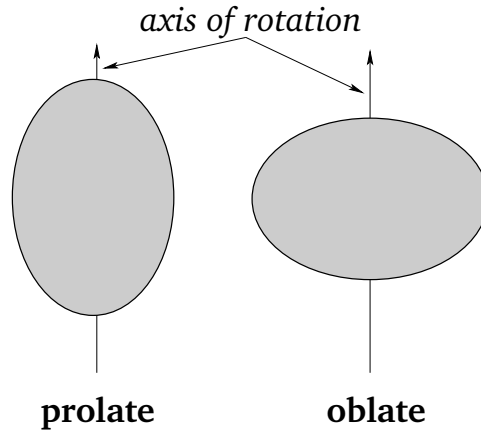


Figure 1.18: Prolate and oblate spheroids.

R . In other words, the slight squashing or elongation of the spheroid along its axis, as ϵ is varied, does not modify its volume.

Let $\Phi(r, \theta)$ and $\rho(r, \theta)$ be the gravitational potential and the mass density of the spheroid, respectively. Let us write

$$\Phi(r, \theta) = \Phi_0(r) + \Phi_2(r) P_2(\cos \theta), \quad (1.331)$$

and

$$\rho(r, \theta) = \rho_0(r) + \rho_2(r) P_2(\cos \theta), \quad (1.332)$$

where

$$\rho_0(r) = \begin{cases} \gamma & r \leq R \\ 0 & r > R \end{cases}, \quad (1.333)$$

and

$$\rho_2(r) = -\frac{2}{3} \gamma R \epsilon \delta(r - R). \quad (1.334)$$

Here, $\delta(x)$ is a *Dirac delta function*. (See Section 2.1.6.) This function has the unusual property that $\delta(x) = 0$ for $x \neq 0$, $\delta(x) = \infty$ at $x = 0$, but

$$\int_{-\infty}^{\infty} \delta(x) dx = 1. \quad (1.335)$$

Thus, a Dirac delta function is an integrable spike function, centered on $x = 0$, that has unit area under it. Note that $\rho_0(r)$ is the density distribution of a uniform sphere of density γ and radius R . On the other hand, $\rho_2(r) P_2(\cos \theta) = \gamma [R_\theta(\theta) - R] \delta(r - R)$ is the density distribution obtained by taking the slight excess or deficit of surface mass, due to the deviation from sphericity of the spheroid, and placing it all at radius R . Note that, in writing Equation (1.331), we have assumed that an axisymmetric mass distribution (i.e., a distribution that is independent of the azimuthal angle, ϕ) gives rise to an axisymmetric gravitational potential.

Now, in spherical coordinates, the Laplacian of $\Phi(r, \theta)$ (i.e., a function of spherical coordinates that is independent of the azimuthal angle, ϕ) takes the form

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial \Phi}{\partial r} \right) + \frac{1}{r^2} \frac{\partial}{\partial \mu} \left[(1 - \mu^2) \frac{\partial \Phi}{\partial \mu} \right], \quad (1.336)$$

where $\mu = \cos \theta$. (See Section A.23.) Thus, according to Poisson's equation, (1.273), we can write

$$\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Phi_0}{dr} \right) = 4\pi G \rho_0, \quad (1.337)$$

$$\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Phi_2}{dr} \right) - 6 \Phi_2 = 4\pi G \rho_2. \quad (1.338)$$

Here, we have made use of the easily proved result

$$\frac{d}{d\mu} \left[(1 - \mu^2) \frac{dP_2(\mu)}{d\mu} \right] = -6 P_2(\mu). \quad (1.339)$$

We have also employed the readily demonstrated result that

$$\int_{-1}^1 P_2(\mu) d\mu = 0, \quad (1.340)$$

which allows us to separately equate the components of Poisson's equation that are independent of θ , and that vary with θ as $P_2(\cos \theta)$.

Equations (1.333) and (1.337) must be solved subject to the physical boundary conditions that $\Phi_0(r)$ be finite at both $r = 0$ and $r = \infty$, and that $\Phi_0(r)$ and its first derivative both be continuous at $r = R$. The latter constraint ensures that the gravitational acceleration is both finite and continuous at $r = R$. It is easily seen, by inspection, that the appropriate solution is

$$\Phi_0(r) = \frac{GM}{2R} \left[\left(\frac{r}{R} \right)^2 - 3 \right] \quad (1.341)$$

for $r < R$, and

$$\Phi_0(r) = -\frac{GM}{r} \quad (1.342)$$

for $r > R$. Here, $M = (4\pi/3)\gamma R^3$ is the mass of the spheroid. If we calculate the associated gravitational acceleration, $\mathbf{g} = -\nabla\Phi_0$, then we can see that this solution is the same as that found for a uniform sphere in Section 1.8.3 using Gauss's law.

Equation (1.334) and (1.338) yield

$$\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Phi_2}{dr} \right) - 6 \Phi_2 = -\frac{8\pi}{3} G \gamma R \epsilon \delta(r - R). \quad (1.343)$$

Now, $\Phi_2(r)$ must be continuous across $r = R$ otherwise the gravitational acceleration would be infinite. Hence, integrating the previous equation across $r = R$, and making use of Equation (1.335), we obtain

$$\left[\frac{d\Phi_2}{dr} \right]_{r=R-}^{r=R+} = -\frac{8\pi}{3} G \gamma R \epsilon. \quad (1.344)$$

Note that the discontinuity in the gradient of Φ_2 at $r = R$ is just an artifact of the fact that we have placed all of the excess or deficit of surface mass at this radius, and is not a real phenomenon (i.e., the discontinuity would be resolved if we were to spread the mass out slightly). Now, for $r \neq R$, the Dirac delta function, $\delta(r - R)$, is zero, and Equation (1.343) reduces to

$$\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Phi_2}{dr} \right) - 6 \Phi_2 = 0. \quad (1.345)$$

This equation must be solved subject to the physical boundary conditions that $\Phi_2(r)$ be finite at both $r = 0$ and $r = \infty$, that $\Phi_2(r)$ be continuous at $r = R$, and that $\Phi_2(r)$ satisfy Equation (1.344). It can be seen, by inspection, that the appropriate solution is

$$\Phi_2(r) = \frac{2}{5} \epsilon \frac{G M r^2}{R^3} \quad (1.346)$$

for $r < R$, and

$$\Phi_2(r) = \frac{2}{5} \epsilon \frac{G M R^2}{r^3} \quad (1.347)$$

for $r > R$.

Hence, we deduce that the net gravitational potential is

$$\Phi(r, \theta) = \frac{G M}{2 R} \left[\left(\frac{r}{R} \right)^2 - 3 \right] + \frac{2}{5} \epsilon \frac{G M r^2}{R^3} P_2(\cos \theta) \quad (1.348)$$

inside the spheroid, and

$$\Phi(r, \theta) = -\frac{G M}{r} + \frac{2}{5} \epsilon \frac{G M R^2}{r^3} P_2(\cos \theta) \quad (1.349)$$

outside the spheroid. In particular, the gravitational potential on the surface of the spheroid is

$$\Phi(R, \theta) = -\frac{G M}{R} \left[1 + \frac{4}{15} \epsilon P_2(\cos \theta) + \mathcal{O}(\epsilon^2) \right]. \quad (1.350)$$

1.10.2 Rotational Flattening of Earth

The Earth rotates diurnally with an angular velocity vector, $\boldsymbol{\Omega}$, that is directed from the center of the Earth toward its north geographic pole, and is of magnitude

$$\Omega = \frac{2\pi}{23^{\text{h}} 56^{\text{m}} 04^{\text{s}}} = 7.292 \times 10^{-5} \text{ rad s}^{-1}. \quad (1.351)$$

Here, $23^{\text{h}} 56^{\text{m}} 04^{\text{s}}$ is the length of a so-called *sidereal day*, and is the period of the Earth's diurnal rotation relative to the distant stars (as opposed to the Sun).

Let the Earth's axis of rotation correspond to the z -axis, and let us set up a conventional spherical coordinate system whose origin is the Earth's center, and whose symmetry axis corresponds to the z -axis. (See Section A.23.) A general point in the Earth whose spherical coordinates are r , θ , ϕ rotates at angular velocity $\Omega \mathbf{e}_z$ in a circle of radius $r \sin \theta$. See Figure 1.19. Thus, according to

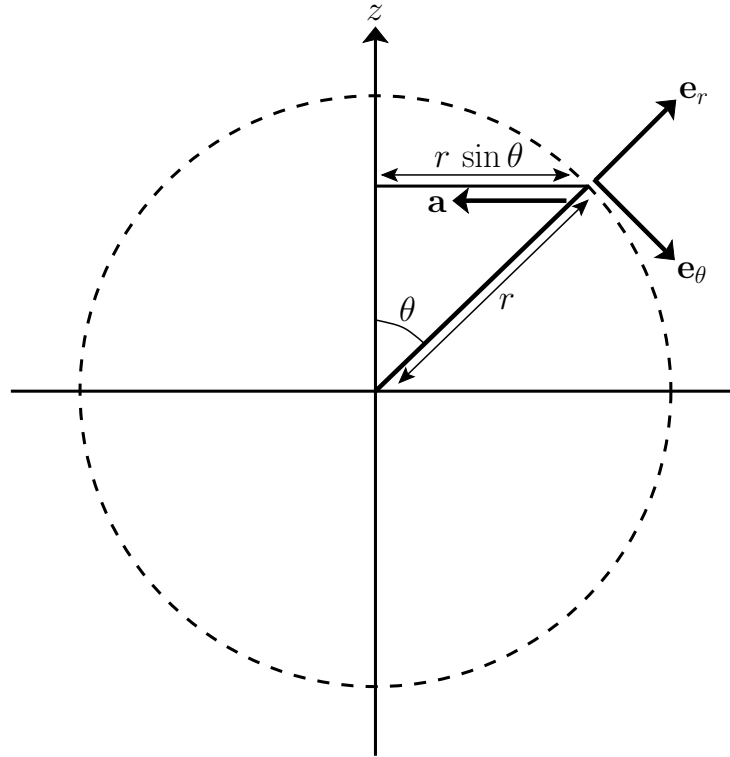


Figure 1.19: Centripetal acceleration.

elementary physics, the point accelerates toward the z -axis with an acceleration $a = \Omega^2 r \sin \theta$. As is clear from Figure 1.19, the point's vector acceleration is

$$\mathbf{a} = -\Omega^2 r \sin \theta (\sin \theta \mathbf{e}_r + \cos \theta \mathbf{e}_\theta), \quad (1.352)$$

where $\mathbf{e}_r = \nabla r / |\nabla r|$ and $\mathbf{e}_\theta = \nabla \theta / |\nabla \theta|$ are unit vectors in the spherical coordinate system. However, it is easily demonstrated that

$$\mathbf{a} = \nabla \chi \quad (1.353)$$

(see Section A.23), where

$$\chi(r, \theta) = -\frac{\Omega^2 r^2}{2} \sin^2 \theta = \frac{\Omega^2 r^2}{3} [P_2(\cos \theta) - 1] \quad (1.354)$$

can be thought of as a kind of centrifugal potential. (Because in the non-inertial frame of reference that co-rotates with the Earth the point in question would appear to be subject to a fictitious centrifugal force $-\nabla \chi$.)

Let us model the interior of the Earth as a fluid of uniform mass density γ . [It turns out that the centrifugal potential, (1.354), is sufficiently large that the rigidity of the rock that makes up the Earth is insufficient to prevent the Earth from responding to the potential in a fluid-like manner.] Now, if $p(\mathbf{r})$ is the pressure distribution in the interior of the Earth then a small cuboid volume of

the Earth lying between x and $x + dx$, y and $y + dy$, and z and $z + dz$, experiences a net pressure force

$$\begin{aligned} \mathbf{F} &= [p(x, y, z) - p(x + dx, y, z)] dy dz \mathbf{e}_x + [p(x, y, z) - p(x, y + dy, z)] dx dz \mathbf{e}_y \\ &\quad + [p(x, y, z) - p(x, y, z + dz)] dx dy \mathbf{e}_z \\ &= -\frac{\partial p}{\partial x} dV \mathbf{e}_x - \frac{\partial p}{\partial y} dV \mathbf{e}_y - \frac{\partial p}{\partial z} dV \mathbf{e}_z \\ &= -\nabla p dV, \end{aligned} \quad (1.355)$$

where $dV = dx dy dz$ is the volume of the cuboid. (See Section A.19). Thus, the force per unit mass due to the pressure inside the Earth is

$$\mathbf{f} = \frac{\mathbf{F}}{\gamma dV} = -\frac{\nabla p}{\gamma}. \quad (1.356)$$

The equation of motion of a general point inside the Earth is

$$\mathbf{a} = -\nabla \Phi - \frac{\nabla p}{\gamma}, \quad (1.357)$$

Here, the first term on the right-hand side of the previous equation is the gravitational force per unit mass acting at the point [see Equation (1.269)], whereas the second term is the force per unit mass due to internal pressure. Making use of Equation (1.353), we deduce that force balance inside the Earth requires that

$$\nabla \left(\Phi + \chi + \frac{p}{\gamma} \right) = 0. \quad (1.358)$$

The previous equation can be integrated to give

$$\Phi + \chi + \frac{p}{\gamma} = c, \quad (1.359)$$

where c is a constant.

Let us model the Earth as a spheroid whose outer radius, $R_\theta(\theta)$, is specified by Equation (1.328). (See Section 1.10.1.) It follows from Equation (1.350) that the gravitational potential at the Earth's surface is

$$\Phi(R_\theta, \theta) \simeq -\frac{GM}{R} \left[1 + \frac{4}{15} \epsilon P_2(\cos \theta) \right], \quad (1.360)$$

where M is the Earth's mass, R its mean radius, and ϵ its ellipticity. It is clear from Equation (1.354) that the centrifugal potential at the Earth's surface is

$$\chi(R_\theta, \theta) \simeq \frac{\Omega^2 R^2}{3} [P_2(\cos \theta) - 1]. \quad (1.361)$$

Note that we have neglected the slight difference between R_θ and R , when evaluating the previous expression, because the centrifugal potential is relatively small compared to the gravitational

potential [see Equation (1.366)], and ϵ is also assumed to be small [see Equation (1.367)]. Now, the pressure at the Earth's surface must be zero, otherwise the surface would not be in equilibrium with outer space. (Here, we are neglecting the relatively small pressure due to the atmosphere.) It follows from the previous three equations that, on the surface of the Earth,

$$-\frac{GM}{R} \left[1 + \frac{4}{15} \epsilon P_2(\cos \theta) \right] + \frac{\Omega^2 R^2}{3} [P_2(\cos \theta) - 1] = c. \quad (1.362)$$

We can separately equate the components of the previous equation that are independent of θ , and that vary with θ as $P_2(\cos \theta)$, to give

$$c = -\frac{GM}{R} - \frac{\Omega^2 R^2}{3}, \quad (1.363)$$

and

$$\epsilon = \frac{15}{4} \zeta. \quad (1.364)$$

where

$$\zeta = \frac{\Omega^2 R^3}{3GM} \quad (1.365)$$

is the ratio of the typical centrifugal acceleration to the typical gravitational acceleration at $r = R$.

Given that $\Omega = 7.292 \times 10^{-5} \text{ rad s}^{-1}$, $R = 6.371 \times 10^6 \text{ m}$, and $M = 5.972 \times 10^{24} \text{ kg}$, we deduce that

$$\zeta = 1.15 \times 10^{-3}, \quad (1.366)$$

and

$$\epsilon = 4.31 \times 10^{-3}. \quad (1.367)$$

In other words, as consequence of the Earth's rotation, the shape of the Earth is an oblate (because $\epsilon > 0$) spheroid. Thus, the Earth is slightly flattened along an axis passing through its geographic poles. This result was first obtained by Newton. The predicted difference between the Earth's equatorial and polar radii is $\Delta R = R_e - R_p = \epsilon R = 27.5 \text{ km}$. In fact, the observed ellipticity of the Earth is

$$\epsilon = 3.35 \times 10^{-3}, \quad (1.368)$$

with an associated difference between the equatorial and polar radii of 21.4 km. Our analysis has overestimated the Earth's rotational flattening because, for the sake of simplicity, we modeled the Earth as a uniform body. In fact, the interior of the Earth is much denser than its crust.

1.10.3 Surface Gravity of Earth

Making use of Equations (1.329), (1.349), (1.354), (1.364), and (1.365), the combined gravitational and centrifugal potential outside the Earth is

$$\Phi + \chi = -\frac{GM}{R} \left[\frac{R}{r} - \frac{3\zeta}{2} \left(\frac{R}{r} \right)^3 + \frac{9\zeta}{4} \left(\frac{R}{r} \right)^3 \sin^2 \theta + \frac{3\zeta}{2} \left(\frac{r}{R} \right)^2 \sin^2 \theta \right]. \quad (1.369)$$

According to Equations (1.328), (1.329), and (1.364), the surface of the Earth lies at radius

$$R_\theta(\theta) = R \left(1 - \frac{5\zeta}{2} + \frac{15\zeta}{4} \sin^2 \theta \right). \quad (1.370)$$

The effective gravitational acceleration at the Earth's surface is $g = |\nabla(\Phi + \chi)|(R_\theta, \theta)$, which reduces to

$$g(\lambda) = \frac{GM}{R^2} \left[1 + \frac{\zeta}{2} - \frac{15\zeta}{4} \cos^2 \lambda + \mathcal{O}(\zeta^2) \right], \quad (1.371)$$

where $\lambda = \pi/2 - \theta$ corresponds to terrestrial latitude. The previous expression shows that the Earth's rotation, combined with its equatorial bulge, causes the acceleration experienced by objects close to the Earth's surface, that co-rotate with the Earth, to vary slightly with latitude. The acceleration is greatest at the poles ($\lambda = \pi/2$), and weakest at the equator ($\lambda = 0$). To be more exact, we predict that

$$g_{\text{pole}} = \frac{GM}{R^2} \left(1 + \frac{\zeta}{2} \right) = 9.826 \text{ m s}^{-1}, \quad (1.372)$$

and

$$g_{\text{equator}} = \frac{GM}{R^2} \left(1 - \frac{13\zeta}{4} \right) = 9.783 \text{ m s}^{-1}. \quad (1.373)$$

In fact, the true values of the polar and equatorial accelerations are $g_{\text{pole}} = 9.832 \text{ m s}^{-1}$ and $g_{\text{equator}} = 9.781 \text{ m s}^{-1}$, respectively.

1.10.4 MacCullagh's Formula

Consider the uniform spheroid discussed in Section 1.10.1. Let us set up a Cartesian coordinate system whose origin coincides with the center of the spheroid, and whose z -axis corresponds to the spheroid's rotation axis.

It is clear, by symmetry, that the x -, y -, and z -axes are principal axes of rotation. (See Section 1.7.2.) Hence, given that $x = r \sin \theta \cos \phi$, $y = r \sin \theta \sin \phi$, $z = r \cos \theta$, where r , θ , z are conventional spherical coordinates (see Section A.23), the principal moments of inertia of the spheroid are

$$\begin{aligned} I_{xx} &= \int_V (y^2 + z^2) \rho(x, y, z) dV = \int_V r^2 (\sin^2 \theta \sin^2 \phi + \cos^2 \theta) \rho(r, \theta) dV \\ &= \frac{1}{2} \int_V r^2 (1 + \cos^2 \theta) \rho(r, \theta) dV, \end{aligned} \quad (1.374)$$

$$\begin{aligned} I_{yy} &= \int_V (x^2 + z^2) \rho(x, y, z) dV = \int_V r^2 (\sin^2 \theta \cos^2 \phi + \cos^2 \theta) \rho(r, \theta) dV \\ &= \frac{1}{2} \int_V r^2 (1 + \cos^2 \theta) \rho(r, \theta) dV, \end{aligned} \quad (1.375)$$

$$I_{zz} = \int_V (x^2 + y^2) \rho(x, y, z) dV = \int_V r^2 (1 - \cos^2 \theta) dV, \quad (1.376)$$

where the integral is over the volume, V , of the spheroid. In writing the previous formulae, we have made use of the fact that the mass density of the spheroid, $\rho(r, \theta)$, is independent of the azimuthal angle, ϕ , and that the averages of $\cos^2 \phi$ and $\sin^2 \phi$ over the volume of a mass distribution that is independent of ϕ are both equal to $1/2$.

Let us write $I_{xx} = I_{yy} = I_{\perp}$ and $I_{\parallel} = I_{zz}$. (See Section 1.7.2.) Now, according to Equations (1.332)–(1.334), the mass density of the spheroid is

$$\rho(r, \theta) = \rho_0(r) - \frac{2}{3} \gamma R \epsilon \delta(r - R) P_2(\cos \theta), \quad (1.377)$$

where

$$\rho_0(r) = \begin{cases} \gamma & r \leq R \\ 0 & r > R \end{cases}. \quad (1.378)$$

Here, R is the mean radius of the spheroid, and γ is its uniform mass density. Given that $dV = r^2 \sin \theta dr d\theta d\phi$, we can evaluate the integrals (1.374)–(1.376) to give

$$I_{\parallel} = \frac{2}{5} M R^2 + \frac{4}{15} \epsilon M R^2, \quad (1.379)$$

$$I_{\perp} = \frac{2}{5} M R^2 - \frac{2}{15} \epsilon M R^2. \quad (1.380)$$

Here, $M = (4\pi/3) R^3 \gamma$ is the mass of the spheroid.

If we apply the previous two results to the Earth, for which $\epsilon = 3.35 \times 10^{-3}$ (see Section 1.10.2), then it is clear that the slight rotational flattening of the Earth along its axis of diurnal rotation causes the principal moment of inertia about this axis, I_{\parallel} , to be slightly larger than the principal moment of inertia, I_{\perp} , about an axis that lies in the Earth's equatorial plane. In particular,

$$I_{\parallel} - I_{\perp} = \frac{2}{5} \epsilon M R^2. \quad (1.381)$$

Thus, it follows from Equation (1.349) that the self-generated gravitational potential outside the Earth can be written

$$\Phi(r, \theta) = -\frac{GM}{r} + \frac{G(I_{\parallel} - I_{\perp})}{r^3} P_2(\cos \theta). \quad (1.382)$$

This result is known as *MacCullagh's formula*. It turns out that the previous formula applies to any axisymmetric mass distribution (i.e., it is not limited to spheroidal mass distributions of uniform mass density). The first term on the right-hand side of MacCullagh's formula is the *monopole* gravitational potential that would be generated if all the Earth's mass were concentrated at its center of mass, whereas the second term is the *quadrupole* gravitational potential generated by the slight deviation of the Earth's shape from that of a sphere.

1.10.5 Gravitational Torque on Axisymmetric Mass Distribution

Consider an axisymmetric mass distribution. Let us set up a Cartesian coordinate system whose origin corresponds to the center of mass of the distribution, and whose z -axis corresponds to the

symmetry axis. The fact that the origin corresponds to the center of mass implies that

$$\int_V x\rho dV = \int_V y\rho dV = \int_V z\rho dV = 0, \quad (1.383)$$

where the integrals are over the volume, V , of the distribution, and ρ is the distribution's mass density. (See Section 1.4.2.) By symmetry, the Cartesian axes are principal axes of rotation, which implies that all of the products of inertia are zero. (See Section 1.7.2.) In other words,

$$\int_V xy\rho dV = \int_V yz\rho dV = \int_V zx\rho dV = 0. \quad (1.384)$$

Finally, the principal moments of inertia of the distribution are

$$I_{\parallel} = I_{zz} = \int_V (x^2 + y^2)\rho dV, \quad (1.385)$$

$$I_{\perp} = I_{xx} = I_{yy} = \int_V (x^2 + z^2)\rho dV = \int_V (y^2 + z^2)\rho dV. \quad (1.386)$$

Suppose that the mass distribution lies in the gravitational field of a distant object. Let $\Phi(x, y, z)$ be the gravitational potential generated by the distant object. The gravitational torque, about the center of mass, that the distant object exerts on the mass distribution is

$$\boldsymbol{\tau} = - \int_V \mathbf{r} \times \nabla \Phi \rho dV. \quad (1.387)$$

[See Section 1.4.5 and Equation (1.269).] However, we expect $\Phi(x, y, z)$ to only vary slightly across the mass distribution, assuming that the distance of the distant object from the origin is much larger than the dimensions of the mass distribution. Thus, we can Taylor expand $\Phi(x, y, z)$ about the origin to give

$$\Phi(x, y, z) \simeq \Phi_x x + \Phi_y y + \Phi_z z + \frac{1}{2} \Phi_{xx} x^2 + \frac{1}{2} \Phi_{yy} y^2 + \frac{1}{2} \Phi_{zz} z^2 + \Phi_{xy} xy + \Phi_{yz} yz + \Phi_{zx} zx. \quad (1.388)$$

Here, Φ_x denotes $\partial\Phi/\partial x$ evaluated at the origin, whereas Φ_{xy} denotes $\partial^2\Phi/\partial x \partial y$ evaluated at the origin, et cetera. Note that we have set the value of Φ at the origin to zero, as we are free to do, given that gravitational potential (like gravitational potential energy) is undefined to an arbitrary additive constant. The previous six equations can be combined to give

$$\tau_x = -(I_{\parallel} - I_{\perp}) \Phi_{yz}, \quad (1.389)$$

$$\tau_y = (I_{\parallel} - I_{\perp}) \Phi_{zx}, \quad (1.390)$$

$$\tau_z = 0. \quad (1.391)$$

We conclude that the distant object exerts a torque on the mass distribution whose direction lies in the distribution's equatorial plane. Note, however, that the distant object is incapable of exerting a torque on a spherically symmetric mass distribution (i.e., a distribution for which $I_{\parallel} = I_{\perp}$).

1.10.6 Lunisolar Precession of Earth

Let us investigate the influence of the Sun on the Earth's diurnal rotation. Consider the Earth-Sun system. See Figure 1.20. From a geocentric viewpoint, the Sun orbits the Earth counterclockwise (if we look from the north), once per year, in an approximately circular orbit of radius $a_s = 1.496 \times 10^{11}$ m. In astronomy, the plane of the Sun's apparent orbit relative to the Earth is known as the *ecliptic plane*. Let us define non-rotating (with respect to distant stars) Cartesian coordinates, centered on the Earth, which are such that the x' - and y' -axes lie in the ecliptic plane, and the z' -axis is normal to this plane (in the sense that the Earth's north pole lies at positive z'). It follows that the z' -axis is directed toward a point in the sky (located in the constellation Draco) known as the *north ecliptic pole*. In the following, we shall treat the x' , y' , z' coordinate system as inertial. This is a reasonable approximation because the orbital acceleration of the Earth is much smaller than the acceleration due to its diurnal rotation. It is convenient to parameterize the instantaneous position of the Sun in terms of a counterclockwise (if we look from the north) azimuthal angle λ_s that is zero on the positive x' -axis. See Figure 1.20. Thus, the coordinates of the Sun in the x' , y' , z' system are

$$x'_s = a_s \cos \lambda_s, \quad (1.392)$$

$$y'_s = a_s \sin \lambda_s, \quad (1.393)$$

$$z'_s = 0. \quad (1.394)$$

Note that

$$\lambda_s = \omega_s t, \quad (1.395)$$

where

$$\omega_s = \left(\frac{G M_s}{a_s^3} \right)^{1/2} \quad (1.396)$$

is the Sun's apparent orbital angular velocity about the Earth. [See Equation (1.323).] Here, $M_s = 1.989 \times 10^{30}$ kg is the mass of the Sun.

Let $\mathbf{\Omega}$ be the Earth's angular velocity vector due to its diurnal rotation. This vector subtends an angle θ with the z' -axis, where $\theta = 23.44^\circ$ is the mean inclination of the ecliptic to the Earth's equatorial plane. Suppose that the projection of $\mathbf{\Omega}$ onto the ecliptic plane subtends an angle ϕ with the x' -axis, where ϕ is measured in a counterclockwise (if we look from the north) sense, and is zero on the positive x' -axis. See Figure 1.20. The orientation of the Earth's axis of rotation (which is, of course, parallel to $\mathbf{\Omega}$) is thus determined by the two angles θ and ϕ . The components of $\mathbf{\Omega}$ in the x' , y' , z' system are

$$\Omega_{x'} = \Omega \sin \theta \cos \phi, \quad (1.397)$$

$$\Omega_{y'} = \Omega \sin \theta \sin \phi, \quad (1.398)$$

$$\Omega_{z'} = \Omega \cos \theta. \quad (1.399)$$

Let us define a second coordinate system, centered on the Earth, such that the z -axis corresponds to the Earth's axis of rotation. The transformation between the x , y , z system and the x' , y' ,

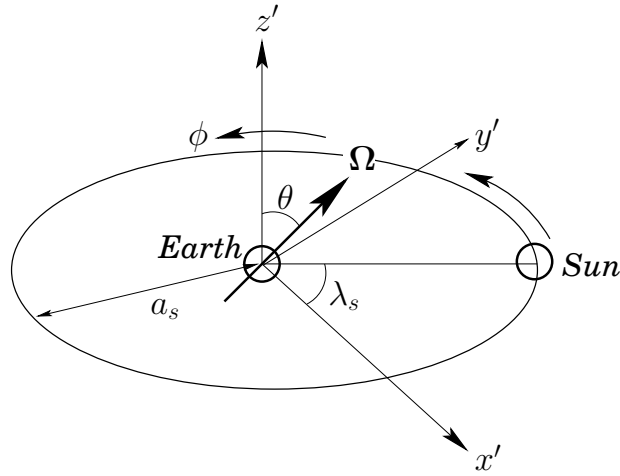


Figure 1.20: The Earth-Sun system.

z' system is

$$\begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} \cos \theta \cos \phi & \cos \theta \sin \phi & -\sin \theta \\ -\sin \phi & \cos \phi & 0 \\ \sin \theta \cos \phi & \sin \theta \sin \phi & \cos \theta \end{pmatrix} \begin{pmatrix} x' \\ y' \\ z' \end{pmatrix}. \quad (1.400)$$

(See Section A.5.) The inverse transformation is

$$\begin{pmatrix} x' \\ y' \\ z' \end{pmatrix} = \begin{pmatrix} \cos \theta \cos \phi & -\sin \phi & \sin \theta \cos \phi \\ \cos \theta \sin \phi & \cos \phi & \sin \theta \sin \phi \\ -\sin \theta & 0 & \cos \theta \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}. \quad (1.401)$$

Of course, the previous two transformations also apply to the components of all vectors in the two coordinate systems. It is easily verified, from Equations (1.397)–(1.400), that $\Omega = \Omega \mathbf{e}_z$ in the x, y, z coordinate system. Note that \mathbf{e}_y lies both in the Earth's equatorial plane and the ecliptic plane.

The gravitational potential of the Sun can be written

$$\Phi(x, y, z) = -\frac{G M_s}{|\mathbf{r} - \mathbf{r}_s|} = -\frac{G M_s}{[(x - x_s)^2 + (y - y_s)^2 + (z - z_s)^2]^{1/2}}, \quad (1.402)$$

where

$$x_s = a_s (\cos \theta \cos \phi \cos \lambda_s + \cos \theta \sin \phi \sin \lambda_s), \quad (1.403)$$

$$y_s = a_s (-\sin \phi \cos \lambda_s + \cos \phi \sin \lambda_s), \quad (1.404)$$

$$z_s = a_s (\sin \theta \cos \phi \cos \lambda_s + \sin \theta \sin \phi \sin \lambda_s) \quad (1.405)$$

are the coordinates of the Sun in the x, y, z system. [See Equations (1.392)–(1.394) and Equation (1.400).] It follows that

$$\Phi_{yz} \equiv \left. \frac{\partial^2 \Phi}{\partial y \partial z} \right|_{x=y=z=0} = -\frac{3 G M_s y_s z_s}{a_s^5} = \frac{3 G M_s}{2 a_s^3} \sin \theta \sin[2(\phi - \lambda_s)], \quad (1.406)$$

$$\Phi_{zx} \equiv \left. \frac{\partial^2 \Phi}{\partial z \partial x} \right|_{x=y=z=0} = -\frac{3 G M_s z_s x_s}{a_s^5} = -\frac{3 G M_s}{2 a_s^3} \cos \theta \sin \theta (1 + \cos[2(\phi - \lambda_s)]). \quad (1.407)$$

Because we are primarily interested in the motion of the Earth's axis of rotation on timescales that are much longer than a year, we can average the preceding expressions over the Sun's orbit to give

$$\Phi_{yz} = 0, \quad (1.408)$$

$$\Phi_{zx} = -\frac{3 G M_s}{2 a_s^3} \cos \theta \sin \theta. \quad (1.409)$$

(This follows because the averages of $\cos[2(\phi - \lambda_s)]$ and $\sin[2(\phi - \lambda_s)]$ over a year are both zero.) Equations (1.389)–(1.391), combined with the previous two equations, reveal that the components of the gravitational torque exerted by the Sun on the Earth in the x, y, z system are

$$\tau_x = 0, \quad (1.410)$$

$$\tau_y = -\frac{3}{2} \frac{G M_s (I_{\parallel} - I_{\perp})}{a_s^3} \cos \theta \sin \theta, \quad (1.411)$$

$$\tau_z = 0, \quad (1.412)$$

where I_{\parallel} and I_{\perp} are the Earth's principal moments of inertia. [See Equations (1.385) and (1.386).] Making use of Equation (1.401), the components of the torque in the x', y', z' system are thus

$$\tau_{x'} = \frac{3}{2} \frac{G M_s (I_{\parallel} - I_{\perp})}{a_s^3} \cos \theta \sin \theta \sin \phi, \quad (1.413)$$

$$\tau_{y'} = -\frac{3}{2} \frac{G M_s (I_{\parallel} - I_{\perp})}{a_s^3} \cos \theta \sin \theta \cos \phi, \quad (1.414)$$

$$\tau_{z'} = 0. \quad (1.415)$$

Because the Earth is rotating about the principal axis of rotation whose principal moment of inertia is I_{\parallel} , the Earth's angular momentum can be written $\mathbf{L} = I_{\parallel} \boldsymbol{\Omega}$. (See Section 1.7.2.) Thus, the components of \mathbf{L} in the x', y', z' system are

$$L_{x'} = I_{\parallel} \Omega \sin \theta \cos \phi, \quad (1.416)$$

$$L_{y'} = I_{\parallel} \Omega \sin \theta \sin \phi, \quad (1.417)$$

$$L_{z'} = I_{\parallel} \Omega \cos \theta. \quad (1.418)$$

[See Equations (1.397)–(1.399).]

The angular equation of motion of the Earth is

$$\frac{d\mathbf{L}}{dt} = \boldsymbol{\tau}. \quad (1.419)$$

(See Section 1.7.1.) However, we must solve this equation in the inertial x', y', z' coordinate system, rather than the x, y, z coordinate system. The reason for this is that the Earth's angular

velocity, Ω , rotates about the z' -axis under the action of the gravitational torque exerted by the Sun on the Earth. Hence, the x, y, z coordinate system, which co-moves with the Earth, accelerates with respect to the x', y', z' system. It follows that the x, y, z system is non-inertial. (See Section 1.5.4.) Making use of Equations (1.413)–(1.418), the components of the previous equation in the x', y', z' system are

$$\frac{d}{dt}(I_{\parallel}\Omega \sin\theta \cos\phi) = \frac{3}{2} \frac{GM_s(I_{\parallel} - I_{\perp})}{a_s^3} \cos\theta \sin\theta \sin\phi, \quad (1.420)$$

$$\frac{d}{dt}(I_{\parallel}\Omega \sin\theta \sin\phi) = -\frac{3}{2} \frac{GM_s(I_{\parallel} - I_{\perp})}{a_s^3} \cos\theta \sin\theta \cos\phi, \quad (1.421)$$

$$\frac{d}{dt}(I_{\parallel}\Omega \cos\theta) = 0. \quad (1.422)$$

If we assume Ω and θ are constants, but that ϕ varies in time, then we can solve the previous three equations to give

$$\frac{d\phi}{dt} = -\Omega_{\phi}, \quad (1.423)$$

where

$$\Omega_{\phi} = \frac{3}{2} \frac{GM_s}{\Omega a_s^3} \left(\frac{I_{\parallel} - I_{\perp}}{I_{\parallel}} \right) \cos\theta = \frac{3}{2} \frac{\omega_s^2}{\Omega} \epsilon \cos\theta, \quad (1.424)$$

and use has been made of Equations (1.379), (1.381), and (1.396).

According to Equation (1.423), the gravitational torque exerted by the Sun on the Earth, due to the Earth's slight oblateness, causes the Earth's axis of rotation to precess steadily about the normal to the ecliptic plane at the rate $-\Omega_{\phi}$. This precession is analogous to that of a spinning top discussed in Section 1.7.7. The fact that $-\Omega_{\phi}$ is negative implies that the precession is in the opposite sense to that of the Earth's diurnal rotation and the Sun's apparent orbit about the Earth. The precession period in (sidereal) years is given by

$$T_{\phi}(\text{yr}) = \frac{\omega_s}{\Omega_{\phi}} = \frac{2T_s(\text{day})}{3\epsilon \cos\theta}, \quad (1.425)$$

where $T_s(\text{day}) = \Omega/\omega_s = 366.26$ is the length of a sidereal year in sidereal days. (Sidereal means measured with respect to distant stars.) Thus, given that $\epsilon = 3.35 \times 10^{-3}$ and $\theta = 23.44^\circ$, we obtain

$$T_{\phi} \simeq 79,400 \text{ years}. \quad (1.426)$$

Unfortunately, the observed precession period of the Earth's axis of rotation about the normal to the ecliptic plane is approximately 25,800 years, so something is clearly missing from our model. It turns out that the missing factor is the influence of the Moon.

From a geocentric viewpoint, the Moon orbits the Earth counterclockwise (if we look from the north), once per month, in an approximately circular orbit of radius $a_m = 3.844 \times 10^8$ m. This orbit is inclined at about 5° to the ecliptic plane. However, in the following, we shall ignore this small inclination, and place the Moon's orbit in the ecliptic plane. Analogous analysis to that employed in the preceding part of this section reveals that the gravitational torque exerted by the Moon on

the Earth (averaged over an month) gives rise to additional contribution to the precession rate Ω_ϕ . In fact, by analogy with Equation (1.424), we expect

$$\Omega_\phi = \frac{3}{2\Omega} \left(\frac{GM_s}{a_s^3} + \frac{GM_m}{a_m^3} \right) \left(\frac{I_\parallel - I_\perp}{I_\parallel} \right) \cos \theta. \quad (1.427)$$

Here, $M_m = 7.324 \times 10^{22}$ kg is the mass of the Moon. Now,

$$\omega_m = \left(\frac{GM_e}{a_m^2} \right)^{1/2} \quad (1.428)$$

is the Moon's orbital angular velocity about the Earth. [See Equation (1.323).] Here, $M_e = 5.972 \times 10^{24}$ kg is the mass of the Earth. Making use of the previous equation, as well as Equations (1.379), (1.381), and (1.396), we obtain

$$\Omega_\phi = \frac{3}{2} \frac{\omega_s^2 + \mu_m \omega_m^2}{\Omega} \cos \theta, \quad (1.429)$$

where $\mu_m = M_m/M_e$.

According to Equations (1.423) and (1.429), the combined gravitational torque exerted by the Sun and the Moon on the Earth, due to the Earth's slight oblateness, causes the Earth's axis of rotation to precess steadily about the normal to the ecliptic plane at the rate $-\Omega_\phi$. As before, the negative sign indicates that the precession is in the opposite direction to the (apparent) orbital motion of the Sun and the Moon. The period of this so-called *lunisolar precession* in (sidereal) years is given by

$$T_\phi(\text{yr}) = \frac{\omega_s}{\Omega_\phi} = \frac{2T_s(\text{day})}{3\epsilon(1 + \mu_m/[T_m(\text{yr})]^2) \cos \theta}, \quad (1.430)$$

where $T_m(\text{yr}) = \omega_s/\omega_m = 0.00748$ is the Moon's (sidereal) orbital period in years. Given that $\epsilon = 3.35 \times 10^{-3}$, $\theta = 23.44^\circ$, $T_s(\text{day}) = 366.26$, and $\mu_m = 0.0123$, we obtain

$$T_\phi \simeq 24,800 \text{ years}. \quad (1.431)$$

This prediction is fairly close to the observed precession period of 25,800 years. The main reason that our estimate is slightly inaccurate is because we have neglected to take into account the small eccentricities of the Earth's orbit around the Sun (see Table 1.4) and the Moon's orbit around the Earth, as well as the small inclination of the Moon's orbit to the Earth's.

The point in the sky toward which the Earth's axis of rotation is directed is known as the *north celestial pole*. Currently, this point lies within about a degree of the fairly bright star Polaris, which is consequently sometimes known as the *north star* or *pole star*. See Figure 1.21. It follows that Polaris appears to be almost stationary in the sky, always lying due north, and can thus be used for navigational purposes. Indeed, mariners have relied on the north star for many hundreds of years to determine direction at sea. Unfortunately, because of the precession of the Earth's axis of rotation, the north celestial pole is not a fixed point in the sky, but instead traces out a circle, of angular radius 23.44° , about the north ecliptic pole, with a period of 25,800 years. See Figure 1.21. Hence,

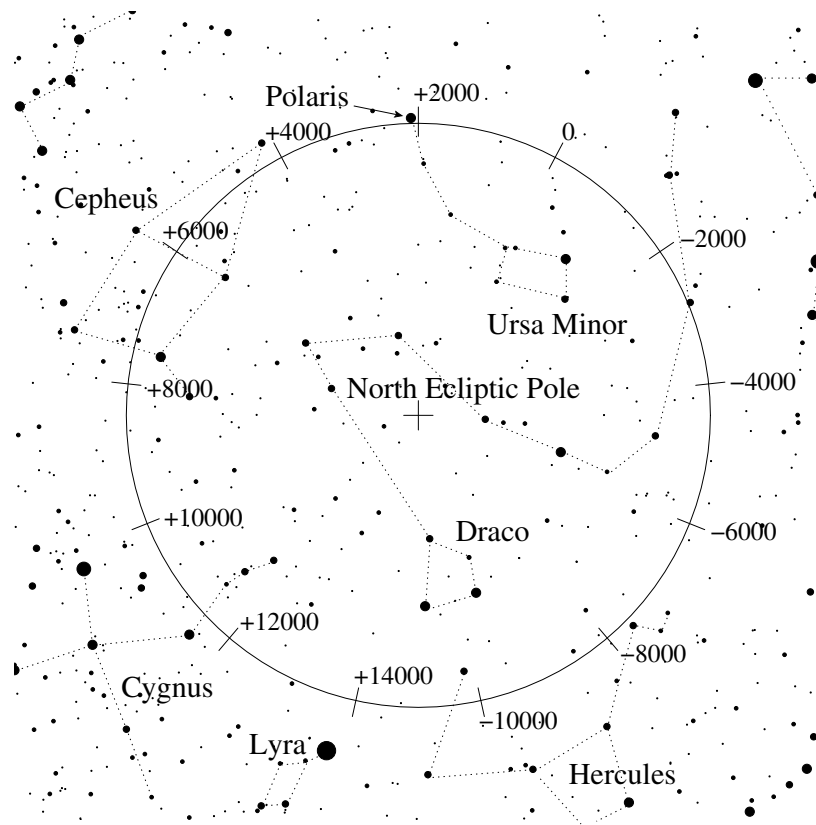


Figure 1.21: Path of the north celestial pole against the backdrop of the stars as consequence of the precession of the equinoxes (calculated assuming constant precessional speed and obliquity). Numbers indicate years relative to start of common era.

a few thousand years from now, the north celestial pole will no longer coincide with Polaris, and there will be no convenient way of telling direction from the stars.

The projection of the ecliptic plane onto the sky is called the *ecliptic circle*, and coincides with the apparent path of the Sun against the backdrop of the stars. The projection of the Earth's equator onto the sky is known as the *celestial equator*. As has been previously mentioned, the ecliptic is inclined at 23.44° to the celestial equator. The two points in the sky at which the ecliptic crosses the celestial equator are called the equinoxes, because night and day are equally long when the Sun lies at these points. Thus, the Sun reaches the vernal equinox on about March 20, and this traditionally marks the beginning of spring (in the Earth's northern hemisphere). Likewise, the Sun reaches the autumnal equinox on about September 22, and this traditionally marks the beginning of autumn. (In fact, in our calculation, the unit vector $\mathbf{e}_y = -\sin \phi \mathbf{e}_{x'} + \cos \phi \mathbf{e}_{y'}$ is directed toward the autumnal equinox.) However, the precession of the Earth's axis of rotation causes the celestial equator (which is always normal to this axis) to precess in the sky; it thus also causes the equinoxes to precess along the ecliptic. This effect is known as the *precession of the equinoxes*. The precession is in the opposite direction to the Sun's apparent motion around the ecliptic, and is of magnitude 1.4° per century. Amazingly, this miniscule effect was discovered

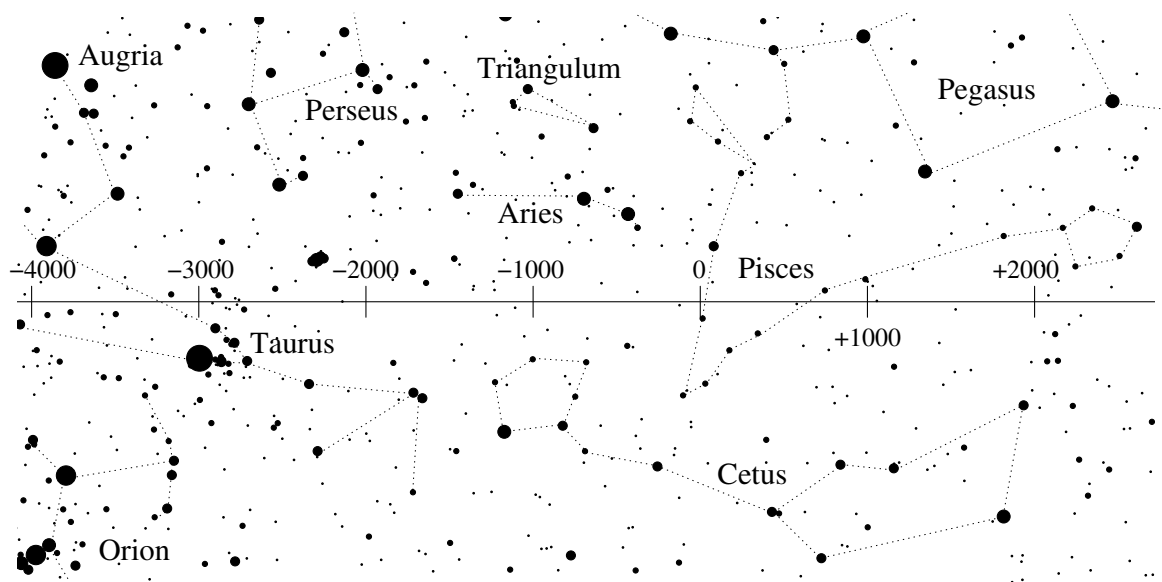


Figure 1.22: Path of the vernal equinox against the backdrop of the stars as a consequence of the precession of the equinoxes (calculated assuming constant precessional speed and obliquity). Numbers indicate years relative to start of common era.

by the ancient Greeks (with the help of ancient Babylonian observations). In about 2000 BCE, when the science of astronomy originated in ancient Egypt and Babylonia, the vernal equinox lay in the constellation Aries. See Figure 1.22. Indeed, the vernal equinox is still sometimes called the *first point of Aries* in astronomical texts. About 90 BCE, the vernal equinox moved into the constellation Pisces, where it still remains. The equinox will move into the constellation Aquarius (marking the beginning of the much heralded “Age of Aquarius”) in about 2600 CE. Incidentally, the position of the vernal equinox in the sky is of great significance in astronomy, because it is used as the zero of celestial longitude (much as the Greenwich meridian is used as the zero of terrestrial longitude).

1.10.7 Two-Body Dynamics

Let us consider the motion of a dynamical system that consists of two freely moving and mutually interacting point objects. Suppose that our first object is of mass m_1 , and is located at displacement \mathbf{r}_1 . Likewise, our second object is of mass m_2 , and is located at displacement \mathbf{r}_2 . Let the first object exert a force \mathbf{f}_{21} on the second. By Newton’s third law, the second object exerts an equal and opposite force, $\mathbf{f}_{12} = -\mathbf{f}_{21}$, on the first. (See Section 1.2.4.) Suppose that there are no other forces in the problem. The equations of motion of our two objects are thus

$$m_1 \frac{d^2 \mathbf{r}_1}{dt^2} = -\mathbf{f}, \quad (1.432)$$

$$m_2 \frac{d^2 \mathbf{r}_2}{dt^2} = \mathbf{f}, \quad (1.433)$$

where $\mathbf{f} = \mathbf{f}_{21}$.

The center of mass of our system is located at

$$\mathbf{R} = \frac{m_1 \mathbf{r}_1 + m_2 \mathbf{r}_2}{m_1 + m_2}. \quad (1.434)$$

(See Section 1.4.2.) Hence, we can write

$$\mathbf{r}_1 = \mathbf{R} - \frac{m_2}{m_1 + m_2} \mathbf{r}, \quad (1.435)$$

$$\mathbf{r}_2 = \mathbf{R} + \frac{m_1}{m_1 + m_2} \mathbf{r}, \quad (1.436)$$

where $\mathbf{r} = \mathbf{r}_2 - \mathbf{r}_1$. Substituting the previous two equations into Equations (1.432) and (1.433), and making use of the fact that the center of mass of an isolated system does not accelerate (see Section 1.4.2), we find that both equations yield

$$\mu \frac{d^2 \mathbf{r}}{dt^2} = \mathbf{f}, \quad (1.437)$$

where

$$\mu = \frac{m_1 m_2}{m_1 + m_2} \quad (1.438)$$

is called the *reduced mass*. Hence, we have effectively converted our original two-body problem into an equivalent one-body problem. In the equivalent problem, the force \mathbf{f} is the same as that acting on both objects in the original problem (modulo a minus sign). However, the mass, μ , is different, and is less than either of m_1 or m_2 (which is why it is called the “reduced” mass).

1.10.8 Binary Star Systems

Approximately half of the stars in our galaxy are members of so-called *binary star systems*. Such systems consist of two stars orbiting about their common center of mass. The distance separating the stars is always very much less than the distance to the nearest-neighbor star. Hence, a binary star system can be treated as a two-body dynamical system to a very good approximation.

In a binary star system, the gravitational force that the first star exerts on the second is

$$\mathbf{f} = -\frac{G m_1 m_2}{r^3} \mathbf{r}, \quad (1.439)$$

where $\mathbf{r} = \mathbf{r}_2 - \mathbf{r}_1$. [See Equation (1.238).] As we have seen, a two-body system can be reduced to an equivalent one-body system whose equation of motion is of the form (1.437), where $\mu = m_1 m_2 / (m_1 + m_2)$. Hence, in this particular case, we can write

$$\frac{m_1 m_2}{m_1 + m_2} \frac{d^2 \mathbf{r}}{dt^2} = -\frac{G m_1 m_2}{r^3} \mathbf{r}, \quad (1.440)$$

which gives

$$\frac{d^2 \mathbf{r}}{dt^2} = -\frac{G M}{r^3} \mathbf{r}, \quad (1.441)$$

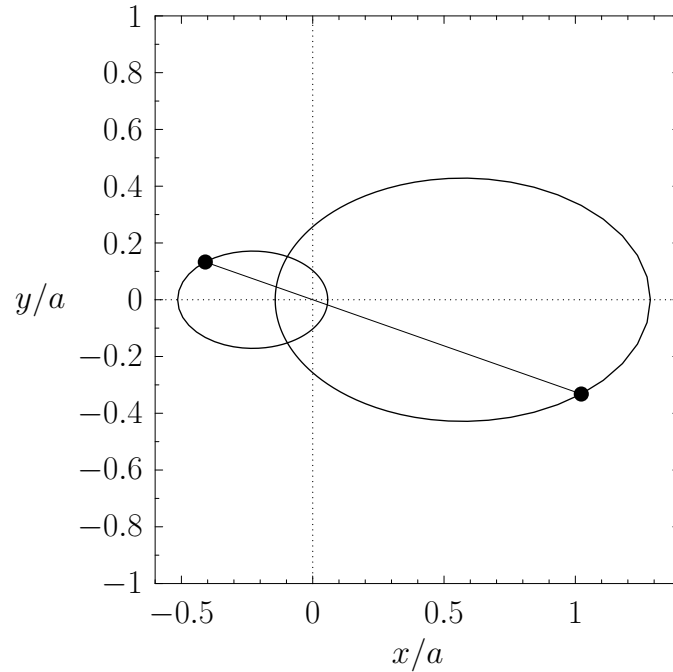


Figure 1.23: An example binary star orbit calculated with $m_1/m_2 = 0.4$ and $e = 0.8$.

where

$$M = m_1 + m_2. \quad (1.442)$$

Equation (1.441) is identical to Equation (1.277), which we have already solved. Hence, we can immediately write down the solution (see Sections 1.9.5–1.9.7):

$$\mathbf{r} = (r \cos \theta, r \sin \theta, 0), \quad (1.443)$$

where

$$r = \frac{a(1 - e^2)}{1 - e \cos \theta}, \quad (1.444)$$

and

$$\frac{d\theta}{dt} = \frac{h}{r^2}, \quad (1.445)$$

with

$$a = \frac{h^2}{(1 - e^2)GM}. \quad (1.446)$$

Here, h is a constant, and we have aligned our Cartesian axes so that the plane of the orbit coincides with the x - y plane. According to the previous solution, the second star executes a Keplerian elliptical orbit, with major radius a and eccentricity e , relative to the first star, and vice versa. From Equation (1.323), the period of revolution, T , is given by

$$T = \sqrt{\frac{4\pi^2 a^3}{GM}}. \quad (1.447)$$

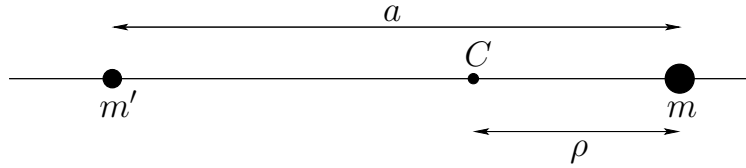


Figure 1.24: Two orbiting masses.

In the inertial frame of reference whose origin always coincides with the center of mass—the so-called *center of mass frame*—the displacement of the two stars are

$$\mathbf{r}_1 = -\frac{m_2}{m_1 + m_2} \mathbf{r}, \quad (1.448)$$

$$\mathbf{r}_2 = \frac{m_1}{m_1 + m_2} \mathbf{r}, \quad (1.449)$$

where \mathbf{r} was specified previously. Figure 1.23 shows an example binary star orbit, in the center of mass frame, calculated with $m_1/m_2 = 0.4$ and $e = 0.8$. It can be seen that both stars execute elliptical orbits about their common center of mass. Furthermore, at any given point in time, the stars are diametrically opposite one another, relative to the center of mass.

Binary star systems have been very useful to astronomers, because it is possible to determine the masses of both stars in such a system by careful observation. The sum of the masses of the two stars, $M = m_1 + m_2$, can be found from Equation (1.447) after a measurement of the major radius, a (which is the mean of the greatest and smallest distance apart of the two stars during their orbit), and the orbital period, T . The ratio of the masses of the two stars, m_1/m_2 , can be determined from Equations (1.448) and (1.449) by observing the fixed ratio of the relative distances of the two stars from the common center of mass about which they both appear to rotate. Obviously, given the sum of the masses, and the ratio of the masses, the individual masses themselves can then be calculated.

1.10.9 Tidal Elongation of Earth

Consider two point objects, of masses m and m' , executing circular orbits about their common center of mass, C , with angular velocity ω . Let a be the distance between the masses, and ρ the distance between point C and mass m . See Figure 1.24. We know from Section 1.10.8 that

$$\omega^2 = \frac{GM}{a^3}, \quad (1.450)$$

and

$$\rho = \frac{m'}{M} a, \quad (1.451)$$

where $M = m + m'$.

Let us transform to a non-inertial frame of reference that rotates, about an axis perpendicular to the orbital plane and passing through C , at the angular velocity ω . In this reference frame, both

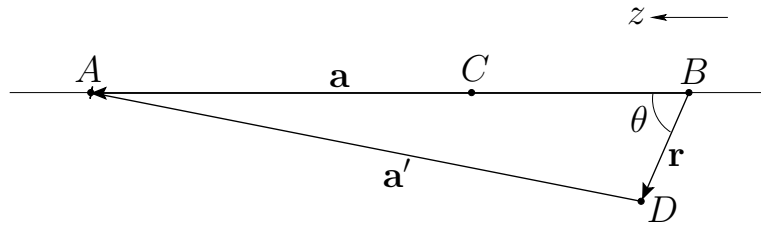


Figure 1.25: Calculation of tidal forces.

masses appear to be stationary. Consider mass m . In the rotating frame, this mass experiences a gravitational acceleration

$$a_g = \frac{G m'}{a^2} \quad (1.452)$$

directed toward the center of mass, and a centrifugal acceleration (see Section 1.10.2)

$$a_c = \omega^2 \rho \quad (1.453)$$

directed away from the center of mass. However, it is easily demonstrated, using Equations (1.450) and (1.451), that

$$a_c = a_g. \quad (1.454)$$

In other words, the gravitational and centrifugal accelerations balance, as must be the case if mass m is to remain stationary in the rotating frame. Let us investigate how this balance is affected if the masses m and m' have finite spatial extents.

Let the center of the mass distribution m' lie at A , the center of the mass distribution m at B , and the center of mass at C . See Figure 1.25. We wish to calculate the centrifugal and gravitational accelerations at some point D in the vicinity of point B . It is convenient to adopt spherical coordinates, centered on point B , and aligned such that the z -axis coincides with the line BA .

Let us assume that the mass distribution m is orbiting around C , but is not rotating about an axis passing through its center of mass, in order to exclude rotational flattening from our analysis. If this is the case then it is easily seen that each constituent point of m executes circular motion of angular velocity ω and radius ρ . See Figure 1.26. Hence, each point experiences the same centrifugal acceleration:

$$\mathbf{g}_c = -\omega^2 \rho \mathbf{e}_z. \quad (1.455)$$

It follows that

$$\mathbf{g}_c = -\nabla \chi', \quad (1.456)$$

where

$$\chi' = \omega^2 \rho z \quad (1.457)$$

is the centrifugal potential and $z = r \cos \theta$. The centrifugal potential can also be written

$$\chi' = \frac{G m'}{a} \frac{r}{a} P_1(\cos \theta), \quad (1.458)$$

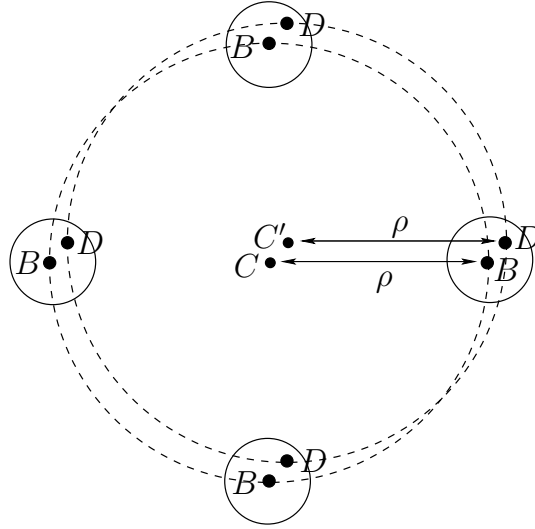


Figure 1.26: The center B of mass distribution m orbits about the center of mass C in a circle of radius ρ . If m is non-rotating then a non-central point D maintains a constant spatial relationship to B , such that D orbits some point C' that has the same spatial relationship to C that D has to B , in a circle of radius ρ .

where

$$P_1(x) = x \quad (1.459)$$

is a Legendre polynomial of degree 1.

The gravitational acceleration at point D due to mass m' is given by

$$\mathbf{g}_g = -\nabla\Phi', \quad (1.460)$$

where the gravitational potential takes the form

$$\Phi' = -\frac{Gm'}{a'}. \quad (1.461)$$

(See Section 1.8.5.) Here, a' is the distance between points A and D . The gravitational potential generated by the mass distribution m' is the same as that generated by an equivalent point mass at A , as long as the distribution is spherically symmetric, which we shall assume to be the case. (See Section 1.8.3.)

Now,

$$\mathbf{a}' = \mathbf{a} - \mathbf{r}, \quad (1.462)$$

where \mathbf{a}' is the vector \overrightarrow{DA} , and \mathbf{a} the vector \overrightarrow{BA} . See Figure 1.25. It follows that

$$a'^{-1} = (a^2 - 2\mathbf{a} \cdot \mathbf{r} + r^2)^{-1/2} = (a^2 - 2ar \cos \theta + r^2)^{-1/2}. \quad (1.463)$$

Expanding in powers of r/a , it is easily demonstrated that

$$\Phi' \simeq -\frac{Gm'}{a} \left[1 + \frac{r}{a} P_1(\cos \theta) + \frac{r^2}{a^2} P_2(\cos \theta) \right], \quad (1.464)$$

to second order in r/a , where the Legendre polynomials $P_2(x)$ and $P_1(x)$ are defined in Equations (1.329) and (1.459), respectively.

Adding χ' and Φ' , we find that

$$\chi = \chi' + \Phi' \simeq -\frac{Gm'}{a} \left[1 + \frac{r^2}{a^2} P_2(\cos \theta) \right], \quad (1.465)$$

to second order in r/a . Note that χ is the potential due to the net externally generated force acting on the mass distribution m in the rotating frame. This potential is constant up to first order in r/a , because the first-order variations in χ' and Φ' cancel each other. The cancellation is a manifestation of the balance between the centrifugal and gravitational accelerations in the equivalent point mass problem discussed previously. However, this balance is only exact at the center of the mass distribution m . Away from the center, the centrifugal acceleration remains constant, whereas the gravitational acceleration increases with increasing z . At positive z , the gravitational acceleration is larger than the centrifugal acceleration, giving rise to a net acceleration in the $+z$ -direction. Likewise, at negative z , the centrifugal acceleration is larger than the gravitational, giving rise to a net acceleration in the $-z$ -direction. It follows that the mass distribution m is subject to a residual acceleration, represented by the second-order variation in Equation (1.465), that acts to elongate it along the z -axis. This effect is known as *tidal elongation*.

Suppose that the mass distribution m is a uniform fluid sphere of radius R . Let us estimate the elongation of this distribution due to the *tidal potential* specified in Equation (1.465), which (neglecting constant terms) can be written

$$\chi(r, \theta) = \frac{Gm}{R} \zeta \left(\frac{r}{R} \right)^2 P_2(\cos \theta). \quad (1.466)$$

Here, the dimensionless parameter

$$\zeta = -\frac{m'}{m} \left(\frac{R}{a} \right)^3 \quad (1.467)$$

is (minus) the typical ratio of the tidal acceleration to the gravitational acceleration at $r \simeq R$. Let us assume that $|\zeta| \ll 1$. By analogy with the analysis in Section 1.10.2, in the presence of the tidal potential, the distribution becomes slightly spheroidal in shape, such that its outer boundary satisfies Equation (1.328). Moreover, the induced ellipticity, ϵ , of the distribution is related to the normalized amplitude, ζ , of the tidal potential according to

$$\epsilon = \frac{15}{4} \zeta. \quad (1.468)$$

[See Equation (1.364).]

Consider the tidal elongation of the Earth due to the Moon. In this case, we have $R = 6.371 \times 10^6$ m, $a = 3.844 \times 10^8$ m, $m = 5.972 \times 10^{24}$ kg, and $m' = 7.324 \times 10^{22}$ kg. Hence, we find that

$$\zeta = -5.58 \times 10^{-8}. \quad (1.469)$$

Thus, according to Equation (1.468), the ellipticity of the Earth induced by the tidal effect of the Moon is

$$\epsilon = \frac{15}{4} \zeta \simeq -2.09 \times 10^{-7}. \quad (1.470)$$

The fact that ϵ is negative implies that the Earth is elongated along the z -axis; that is, along the axis joining its center to that of the Moon. [See Equation (1.328).] If R_+ and R_- are the greatest and least radii of the Earth, respectively, due to this elongation, then

$$\Delta R = R_+ - R_- = -\epsilon R = 1.33 \text{ m.} \quad (1.471)$$

Thus, we predict that the tidal effect of the Moon (which is actually due to spatial gradients in the Moon's gravitational field) causes the Earth to elongate along the axis joining its center to that of the Moon by about 133 centimeters. This turns out to be an overestimate because the tidal potential of the Moon is not strong enough to force the rocks that make up the Earth to respond to it as a fluid.

Consider the tidal elongation of the Earth due to the Sun. In this case, we have $R = 6.371 \times 10^6 \text{ m}$, $a = 1.496 \times 10^{11} \text{ m}$, $m = 5.972 \times 10^{24} \text{ kg}$, and $m' = 1.989 \times 10^{30} \text{ kg}$. Hence, we find that

$$\zeta = -2.65 \times 10^{-8}, \quad (1.472)$$

and

$$\epsilon = \frac{15}{4} \zeta = -9.95 \times 10^{-8}, \quad (1.473)$$

with

$$\Delta R = R_+ - R_- = -\epsilon R = 0.63 \text{ m.} \quad (1.474)$$

Again, this turns out to be an overestimate because the tidal potential of the Sun is not strong enough to force the rocks that make up the Earth to respond to it as a fluid. Nevertheless, we can conclude that the tidal elongation of the Earth due to the Sun is about half that due to the Moon.

Because the Earth's oceans are liquid, their tidal elongation is significantly larger than that of the underlying land. Hence, the oceans rise, relative to the land, in the region of the Earth closest to the Moon, and also in the region furthest away. Because the Earth is rotating, while the tidal bulge of the oceans remains relatively stationary, the Moon's tidal effect causes the ocean at a given point on the Earth's surface to rise and fall twice daily, giving rise to the phenomenon known as the *tides*. There is also an oceanic tidal bulge due to the Sun that is about half as large as that due to the Moon. Consequently, ocean tides are particularly high when the Sun, the Earth, and the Moon lie approximately in a straight line, so that the tidal effects of the Sun and the Moon reinforce one another. This occurs at a new moon, or at a full moon. These type of tides are called *spring tides* (the name has nothing to do with the season). Conversely, ocean tides are particularly low when the Sun, the Earth, and the Moon form a right angle, so that the tidal effects of the Sun and the Moon partially cancel one another. These type of tides are called *neap tides*. Generally speaking, we would expect two spring tides and two neap tides per month.

Chapter 2

Classical Electromagnetism

2.1 Electrostatic Fields

2.1.1 Electricity

We usually associate electricity with the 20th century (CE), during which it revolutionized the lives of countless millions of ordinary people, in much the same manner as steam power revolutionized lives in the 18th century. It is, therefore, somewhat surprising to learn that humans have known about electricity for many thousands of years. In about 1000 BCE, the ancient Greeks started to navigate the Black Sea, and opened up trade routes, via the river Dnieper, to the Baltic region. Amongst the many trade items that the Greeks obtained from the Baltic was a substance that they called “electron” (*ἤλεκτρον*), but that we nowadays call amber. Amber is fossilized pine resin, and was used by the Greeks, much as it is used today, as a gem stone. However, in about 600 BCE, Thales of Miletus discovered that amber possesses a rather peculiar property; namely, when it is rubbed with fur it develops the ability to attract light objects, such as feathers. For many centuries, this strange phenomenon was thought to be a unique property of amber.

In Elizabethan times, the physician William Gilbert coined the word “electric” (from the Greek word for amber) to describe the previously mentioned effect. It was later found that many materials become electric when rubbed with certain other materials. In 1733 (CE), the chemist Charles du Fay discovered that there are, in fact, two different types of electricity. When amber is rubbed with fur it acquires so-called “resinous” electricity. On the other hand, when glass is rubbed with silk it acquires so-called “vitreous” electricity. Electricity repels electricity of the same kind, but attracts electricity of the opposite kind. At the time, it was thought that electricity was created by friction.

Scientists in the 18th century eventually developed the concept of *electric charge* in order to account for a large body of observations made in countless electrical experiments. There are two types of electric charge; positive (which is equivalent to vitreous), and negative (which is equivalent to resinous). Like electric charges repel one another, whereas opposite charges attract. When two bodies are rubbed together, electric charge can be transferred from one to the other, but the total charge remains constant. Thus, when amber is rubbed with fur, there is transfer of electric charge such that the amber acquires a negative charge, and the fur an equal positive charge. Likewise, when glass is rubbed with silk, the glass acquires a positive charge, and the silk an equal negative

charge. The idea that electric charge is a conserved quantity is attributed to the Benjamin Franklin (who is also to blame for the unfortunate sign convention that led to electrons having a negative charge).

In the 20th century, scientists, such as J.J. Thompson and Ernest Rutherford, discovered that the atoms out of which ordinary matter is composed consist of two components; a relatively massive, positively charged nucleus, surrounded by a cloud of relatively light, negatively charged particles called *electrons*. Electrons and atomic nuclei carry fixed electrical charges, and are essentially indestructible (provided that we neglect nuclear reactions). Under normal circumstances, only the electrons are mobile. Thus, when amber is rubbed with fur, electrons are transferred from the fur to the amber, giving the amber an excess of electrons, and, hence, a negative electric charge, and the fur a deficit of electrons, and, hence, a positive charge. Substances normally contain neither an excess nor a deficit of electrons, and are, therefore, electrically neutral.

The SI unit of electric charge is the *coulomb* (C). The electric charge of an electron is

$$e = -1.602 \times 10^{-19} \text{ C.} \quad (2.1)$$

2.1.2 Coulomb's Law

Between 1785 and 1787, Charles Augustine de Coulomb performed a series of experiments involving electric charges, and eventually established what is nowadays known as *Coulomb's law*. According to this law, any two point electric charges (i.e., electrically charged objects of negligible spatial extents) exert a force on one another. This force is directed along the line of centers joining the two charges, is repulsive for two like charges and attractive for opposite charges, is directly proportional to the product of the charges, and is inversely proportional to the square of the distance between the charges.

Consider a system consisting of two point electric charges. Let charge 1 have electric charge q_1 and displacement \mathbf{r}_1 . Let charge 2 have electric charge q_2 and displacement \mathbf{r}_2 . Coulomb's law states that the electrical force exerted on charge 2 by charge 1 is

$$\mathbf{f}_{21} = \frac{q_1 q_2}{4\pi \epsilon_0} \frac{\mathbf{r}_2 - \mathbf{r}_1}{|\mathbf{r}_2 - \mathbf{r}_1|^3}. \quad (2.2)$$

An equal and opposite force acts on the first charge, in accordance with Newton's third law of motion. (See Section 1.2.4.) The universal constant ϵ_0 is called the *electrical permittivity of free space*, and takes the value

$$\epsilon_0 = 8.854 \times 10^{-12} \text{ C}^2 \text{ N}^{-1} \text{ m}^{-2}. \quad (2.3)$$

As we saw in Section 1.8.1, according to Newtonian gravity, if two point mass objects of masses m_1 and m_2 are located at displacements \mathbf{r}_1 and \mathbf{r}_2 , respectively, then the gravitational force acting on the second object is

$$\mathbf{f}_{21} = -G m_1 m_2 \frac{\mathbf{r}_2 - \mathbf{r}_1}{|\mathbf{r}_2 - \mathbf{r}_1|^3}. \quad (2.4)$$

The universal gravitational constant G takes the value

$$G = 6.674 \times 10^{-11} \text{ N m}^2 \text{ kg}^{-2}. \quad (2.5)$$

[See Equation (1.239).] Note that Coulomb's law has the same mathematical form as Newton's law of gravity. In particular, they are both inverse-square force laws; that is,

$$|\mathbf{f}_{21}| \propto \frac{1}{|\mathbf{r}_2 - \mathbf{r}_1|^2}. \quad (2.6)$$

However, Coulomb's and Newton's laws differ in two crucial respects. First, the force due to gravity is always attractive (because there is no such thing as a negative mass). Second, the magnitudes of the forces predicted by the two laws are vastly different. Consider the ratio of the electrical and gravitational forces acting on two particles. This ratio is a constant, independent of the relative positions of the particles, and is given by

$$\frac{|\mathbf{f}_{\text{electrical}}|}{|\mathbf{f}_{\text{gravitational}}|} = \frac{|q_1| |q_2|}{m_1 m_2} \frac{1}{4\pi \epsilon_0 G}. \quad (2.7)$$

For electrons, the charge to mass ratio is $|q|/m = 1.759 \times 10^{11} \text{ C kg}^{-1}$, so

$$\frac{|\mathbf{f}_{\text{electrical}}|}{|\mathbf{f}_{\text{gravitational}}|} = 4.17 \times 10^{42}, \quad (2.8)$$

which is a truly colossal number. Suppose we were studying a physics problem involving the motion of particles under the action of two forces with the same spatial range, but differing in magnitude by a factor 10^{42} . It would seem a plausible approximation (to say the least) to start the investigation by neglecting the weaker force altogether. Applying this reasoning to the motion of particles in the universe, we would expect the universe to be governed entirely by electrical forces. However, this is not the case. The force that holds us to the surface of the Earth, and prevents us from floating off into space, is gravity. The force that causes the Earth to orbit the Sun is also gravity. In fact, on astronomical lengthscales, gravity is the dominant force, and electrical forces are largely irrelevant. The key to understanding this paradox is that there are both positive and negative electric charges, whereas there are only positive gravitational "charges." This implies that gravitational forces are always cumulative, whereas electrical forces can cancel one another out. Suppose, for the sake of argument, that the universe starts out with randomly distributed electric charges. Initially, we expect electrical forces to completely dominate gravitational forces. These forces act to cause every positive electric charge to get as far away as possible from the other positive charges in the universe, and as close as possible to the other negative charges. After a while, we would expect the positive and negative electric charges to form close pairs. Just how close is determined by quantum mechanics, but, in general, it is fairly close; that is, about 10^{-10} m . The electrical forces due to the charges in each pair effectively cancel one another out on lengthscales much larger than the mutual spacing of the pair. However, it is only possible for gravity to be the dominant long-range force in the universe if the number of positive electric charges is almost equal to the number of negative charges. In this situation, every positive charge can find a negative charge to team up with, and there are virtually no charges left over. In order for the cancellation of long-range electrical forces to be effective, the relative difference in the number of positive and negative electric charges in the universe must be incredibly small. In fact, positive and negative charges have to cancel one another to such accuracy that most physicists believe that

the net electric charge of the universe is exactly zero. But, it is not sufficient for the universe to start out with zero net charge. Suppose there were some elementary particle process that did not conserve electric charge. Even if this were to go on at a very low rate, it would not take long before the fine balance between positive and negative charges in the universe was wrecked. Thus, it is important that electric charge is a conserved quantity (i.e., the net charge of the universe can neither increase or decrease). As far as we know, this is the case. To date, no elementary particle reaction has been discovered that can create or destroy net electric charge.

In summary, there are two long-range forces in the universe, electricity and gravity. The former is enormously stronger than the latter, but is usually hidden away inside neutral atoms. The fine balance of forces due to negative and positive electric charges starts to break down on atomic scales. In fact, interatomic and intermolecular forces are all electrical in nature. So, electrical forces are basically what prevent us from falling through the floor. But, this is electromagnetism on the microscopic, or atomic, scale. *Classical electromagnetism* generally describes phenomena in which some sort of violence is done to matter, so that the close pairing of negative and positive electric charges is disrupted, allowing electrical forces to manifest themselves on *macroscopic* lengthscales. Of course, very little disruption is necessary before gigantic forces are generated. Hence, it is no coincidence that the vast majority of useful machines that humankind has devised during the last century or so are electrical in nature.

2.1.3 Electric Field

Consider a system of N point electric charges. Let the i th charge have electric charge q_i and displacement \mathbf{r}_i . As is the case for gravitational forces (see Section 1.8.1), it is an experimentally demonstrated fact that electrical forces are superposable; that is, the electrical force acting on a test charge whose electric charge is q and whose displacement is \mathbf{r} is simply the sum of all of the Coulomb-law forces exerted on it by each of the other N charges taken in isolation. In other words, the electrical force exerted by the i th charge (say) on the test charge is the same as if all of the other charges were not present. Thus, generalizing Equation (2.2), the force acting on the test charge is given by

$$\mathbf{f}(\mathbf{r}) = q \sum_{i=1, N} \frac{q_i}{4\pi \epsilon_0} \frac{\mathbf{r} - \mathbf{r}_i}{|\mathbf{r} - \mathbf{r}_i|^3}. \quad (2.9)$$

It is helpful to introduce a vector field, $\mathbf{E}(\mathbf{r})$, known as the *electric field*, which is defined as the force exerted on a test charge of unit electric charge whose displacement is \mathbf{r} . Thus, from the previous equation, the electrical force on a test charge q whose displacement is \mathbf{r} is written

$$\mathbf{f}(\mathbf{r}) = q \mathbf{E}(\mathbf{r}), \quad (2.10)$$

where the electric field is given by

$$\mathbf{E}(\mathbf{r}) = \sum_{i=1, N} \frac{q_i}{4\pi \epsilon_0} \frac{\mathbf{r} - \mathbf{r}_i}{|\mathbf{r} - \mathbf{r}_i|^3}. \quad (2.11)$$

At this point, we have no reason to believe that the electric field has any real physical existence. It is just a useful device for calculating the electrical force that acts on test charges placed at various locations.

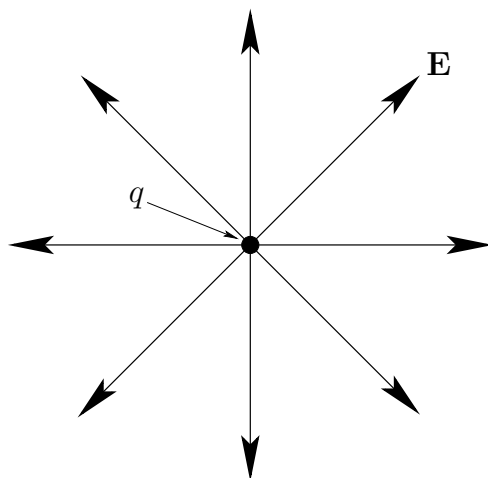


Figure 2.1: Electric field-lines generated by a positive charge.

According to the previous equation, the electric field generated by a single point electric charge q located at the origin is purely radial, is directed outward if the charge is positive, and inward if it is negative, and has magnitude

$$E_r(r) = \frac{q}{4\pi \epsilon_0 r^2}, \quad (2.12)$$

where r is a spherical polar coordinate. Moreover, E_r is the radial component of the field in spherical polar coordinates. The other components are zero. (See Section A.23.) We can represent an electric field by so-called *field-lines*. The direction of the lines indicates the direction of the local electric field, and the density of the lines perpendicular to this direction is proportional to the magnitude of the local electric field. It follows from Equation (2.12) that the number of field-lines crossing the surface of a sphere centered on a point charge (which is equal to E_r times the area, $4\pi r^2$, of the surface) is independent of the radius of the sphere. Thus, the field of a point positive electric charge is represented by a group of equally-spaced, unbroken, straight-lines radiating from the charge. See Figure 2.1. Likewise, field of a point negative charge is represented by a group of equally-spaced, unbroken, straight-lines converging on the charge.

Because electrical forces are superposable, it follows that electric fields are also superposable. In other words, the electric field generated by a collection of electric charges is simply the sum of the fields generated by each of the charges taken in isolation. Suppose that, instead of having a collection of discrete electric charges, we have a continuous distribution of charge represented by an electric charge density $\rho(\mathbf{r})$. Thus, the electric charge at displacement \mathbf{r}' is $\rho(\mathbf{r}') dV'$, where dV' is the volume element at \mathbf{r}' . It follows from a straight-forward extension of Equation (2.11) that the electric field generated by this charge distribution is

$$\mathbf{E}(\mathbf{r}) = \frac{1}{4\pi \epsilon_0} \int_{V'} \rho(\mathbf{r}') \frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3} dV', \quad (2.13)$$

where the volume integral is over a volume, V' , that contains all of the charges.

2.1.4 Electric Scalar Potential

Suppose that $\mathbf{r} = (x, y, z)$ and $\mathbf{r}' = (x', y', z')$ in Cartesian coordinates. The x -component of $(\mathbf{r} - \mathbf{r}')/|\mathbf{r} - \mathbf{r}'|^3$ is written

$$\frac{x - x'}{[(x - x')^2 + (y - y')^2 + (z - z')^2]^{3/2}}. \quad (2.14)$$

However, it is easily demonstrated that

$$\begin{aligned} & \frac{x - x'}{[(x - x')^2 + (y - y')^2 + (z - z')^2]^{3/2}} = \\ & -\frac{\partial}{\partial x} \left(\frac{1}{[(x - x')^2 + (y - y')^2 + (z - z')^2]^{1/2}} \right). \end{aligned} \quad (2.15)$$

Here, $\partial/\partial x$ denotes differentiation with respect to x at constant $y, z, x', y',$ and z' . Because there is nothing special about the x -axis, we can write

$$\frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3} = -\nabla \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right), \quad (2.16)$$

where $\nabla \equiv (\partial/\partial x, \partial/\partial y, \partial/\partial z)$ is a differential operator that involves the components of \mathbf{r} , but not those of \mathbf{r}' . (See Section A.19.) It follows from Equation (2.13) that

$$\mathbf{E} = -\nabla\phi, \quad (2.17)$$

where

$$\phi(\mathbf{r}) = \frac{1}{4\pi\epsilon_0} \int_{V'} \frac{\rho(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} dV'. \quad (2.18)$$

Thus, we conclude that the electric field, $\mathbf{E}(\mathbf{r})$, generated by a collection of fixed electric charges can be written as minus the gradient of a scalar field, $\phi(\mathbf{r})$ —known as the *electric scalar potential*—and that this scalar field can be expressed as a simple volume integral involving the electric charge distribution.

The scalar potential generated by an electric charge q located at the origin is

$$\phi(r) = \frac{q}{4\pi\epsilon_0 r}, \quad (2.19)$$

where r is a spherical polar coordinate. (See Section A.23.) Moreover, according to Equations (2.11) and (2.16), the scalar potential generated by a set of N discrete charges q_i , located at displacements \mathbf{r}_i , is

$$\phi(\mathbf{r}) = \sum_{i=1, N} \phi_i(\mathbf{r}), \quad (2.20)$$

where

$$\phi_i(\mathbf{r}) = \frac{q_i}{4\pi\epsilon_0 |\mathbf{r} - \mathbf{r}_i|}. \quad (2.21)$$

Thus, the net scalar potential is just the sum of the potentials generated by each of the charges taken in isolation.

2.1.5 Electric Potential Energy

Suppose that a particle of electric charge q is taken along some path from point P to point Q . The net work done on the particle by electrical forces is

$$W = \int_P^Q \mathbf{f} \cdot d\mathbf{r}, \quad (2.22)$$

where $\mathbf{f}(\mathbf{r})$ is the electrical force, and $d\mathbf{r}$ is an element of the path. (See Section 1.3.2.) Making use of Equations (2.10) and (2.17), we obtain

$$W = q \int_P^Q \mathbf{E} \cdot d\mathbf{r} = -q \int_P^Q \nabla\phi \cdot d\mathbf{r} = -q [\phi(Q) - \phi(P)]. \quad (2.23)$$

(See Section A.18.) Thus, the work done on the particle is simply minus the product of its charge and the difference in electric potential between the end point and the beginning point. This work is clearly independent of the path taken between points P and Q . Thus, we conclude that an electric field generated by stationary charges is an example of a conservative force field. (See Section 1.3.3.) The work done on the particle when it is taken around a closed loop is zero, so

$$\oint_C \mathbf{E} \cdot d\mathbf{r} = 0 \quad (2.24)$$

for any closed loop C . This implies from the curl theorem that

$$\nabla \times \mathbf{E} = \mathbf{0} \quad (2.25)$$

for any electric field generated by stationary charges. (See Section A.22.) Equation (2.25) also follows directly from Equation (2.17), because $\nabla \times \nabla\phi \equiv \mathbf{0}$ for any scalar potential ϕ . (See Section A.22.)

The SI unit of electric potential is the *volt* (V), which is equivalent to a joule per coulomb. Thus, according to Equation (2.23), the electrical work done on a particle when it is taken between two points is the product of minus its electric charge and the voltage difference between the points.

We are familiar with the idea that a particle moving in a gravitational field possesses potential energy, as well as kinetic energy. (See Section 1.3.5.) If the particle moves from point P to a lower point Q then the gravitational field does work on the particle, causing its kinetic energy to increase. The increase in kinetic energy of the particle is balanced by an equal decrease in its potential energy, so that the overall energy of the particle is a conserved quantity. Therefore, the work done on the particle as it moves from P to Q is minus the difference in its gravitational potential energy between points Q and P . Of course, it only makes sense to talk about gravitational potential energy because the gravitational field is conservative. Thus, the work done in taking a particle between two points is path independent, and, therefore, well defined. This implies that the difference in potential energy of the particle between the beginning and end points is also well defined. We have already seen that an electric field generated by stationary charges is conservative. It follows that we can define an *electric potential energy* of a particle moving in such a field. By

analogy with gravitational fields, the work done in taking a particle of electric charge q from point P to point Q is equal to minus the difference in the electric potential energy of the particle between points Q and P . It follows from Equation (2.23) that the electric potential energy of the particle at a general point Q , relative to some reference point P (where the potential energy is set to zero), is given by

$$W(Q) = q\phi(Q), \quad (2.26)$$

where $\phi(Q)$ is the electric scalar potential at point Q . Free particles tend to move down gradients of potential energy, in order to attain a minimum potential energy state. (See Section 1.3.6.) Thus, free particles in the Earth's gravitational field tend to fall downward. Likewise, positive charges moving in an electric field tend to migrate towards regions with the most negative voltage, and vice versa for negative charges.

The scalar electric potential is undefined to an additive constant. In other words, the transformation

$$\phi(\mathbf{r}) \rightarrow \phi(\mathbf{r}) + c, \quad (2.27)$$

where c is a spatial constant, leaves the electric field unchanged according to Equation (2.17). The scalar potential can be fixed unambiguously by specifying its value at a single point. The usual convention is to say that the potential is zero at infinity. This convention is implicit in Equation (2.18), where it can be seen that $\phi \rightarrow 0$ as $|\mathbf{r}| \rightarrow \infty$, provided that the total electric charge $\int_{V'} \rho(\mathbf{r}') dV'$ is finite.

2.1.6 Gauss's Law

Consider a single electric charge q located at the origin. The electric field generated by such a charge is given by Equation (2.12). Suppose that we surround the charge by a concentric spherical surface S of radius r . See Figure 2.2. The flux of the electric field through this surface is given by

$$\oint_S \mathbf{E} \cdot d\mathbf{S} = \oint_S E_r dS_r = E_r(r) 4\pi r^2 = \frac{q}{4\pi \epsilon_0 r^2} 4\pi r^2 = \frac{q}{\epsilon_0}, \quad (2.28)$$

because the normal to the surface is always parallel to the local electric field. (See Section A.16.) Here, r is also a spherical polar coordinate. (See Section A.23.)

However, we also know from the divergence theorem that

$$\oint_S \mathbf{E} \cdot d\mathbf{S} = \int_V \nabla \cdot \mathbf{E} dV, \quad (2.29)$$

where V is the volume enclosed by surface S . (See Section A.20.) Let us evaluate $\nabla \cdot \mathbf{E}$ directly. In Cartesian coordinates, the electric field (2.12) is written

$$\mathbf{E} = \frac{q}{4\pi \epsilon_0} \left(\frac{x}{r^3}, \frac{y}{r^3}, \frac{z}{r^3} \right), \quad (2.30)$$

where $r^2 = x^2 + y^2 + z^2$. So,

$$\frac{\partial E_x}{\partial x} = \frac{q}{4\pi \epsilon_0} \left(\frac{1}{r^3} - \frac{3x}{r^4} \frac{x}{r} \right) = \frac{q}{4\pi \epsilon_0} \frac{r^2 - 3x^2}{r^5}. \quad (2.31)$$

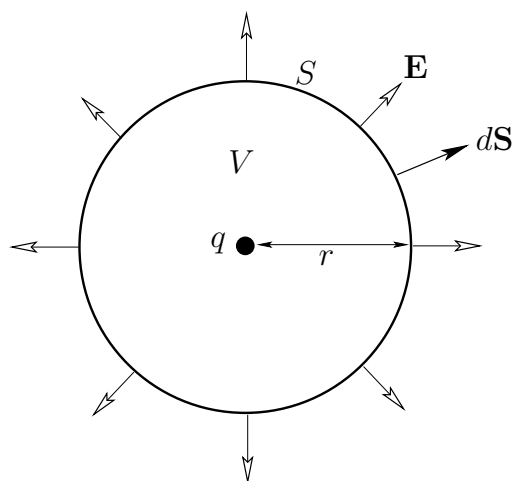


Figure 2.2: Gauss' law.

Here, use has been made of the easily demonstrated result

$$\frac{\partial r}{\partial x} = \frac{x}{r}. \quad (2.32)$$

Formulae analogous to Equation (2.31) can be obtained for $\partial E_y/\partial y$ and $\partial E_z/\partial z$. The divergence of the field is, thus, given by

$$\nabla \cdot \mathbf{E} \equiv \frac{\partial E_x}{\partial x} + \frac{\partial E_y}{\partial y} + \frac{\partial E_z}{\partial z} = \frac{q}{4\pi \epsilon_0} \frac{3r^2 - 3x^2 - 3y^2 - 3z^2}{r^5} = 0. \quad (2.33)$$

(See Section A.20.) This is an extremely puzzling result. We have from Equations (2.28) and (2.29) that

$$\int_V \nabla \cdot \mathbf{E} dV = \frac{q}{\epsilon_0}, \quad (2.34)$$

and yet we have just proved that $\nabla \cdot \mathbf{E} = 0$. This paradox can be resolved after a close examination of Equation (2.33). At the origin ($r = 0$), we find that $\nabla \cdot \mathbf{E} = 0/0$, which implies that $\nabla \cdot \mathbf{E}$ can take any value at this point. Thus, Equations (2.33) and (2.34) can be reconciled if $\nabla \cdot \mathbf{E}$ is some sort of “spike” function; that is, if it is zero everywhere, except arbitrarily close to the origin, where it becomes very large. This must occur in such a manner that the volume integral over the spike is finite.

Let us examine how we might construct a one-dimensional spike function. Consider the “box-car” function

$$g(x, \epsilon) = \begin{cases} 1/\epsilon & \text{for } |x| < \epsilon/2 \\ 0 & \text{otherwise} \end{cases}. \quad (2.35)$$

See Figure 2.3. It is clear that

$$\int_{-\infty}^{\infty} g(x, \epsilon) dx = 1. \quad (2.36)$$

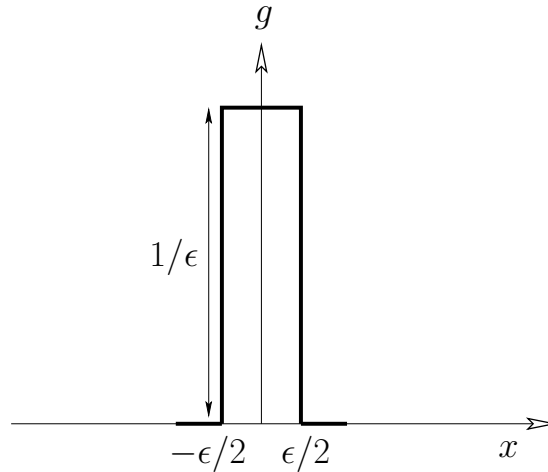


Figure 2.3: A box-car function.

Now, consider the function

$$\delta(x) = \lim_{\epsilon \rightarrow 0} g(x, \epsilon). \quad (2.37)$$

This function is zero everywhere, except arbitrarily close to $x = 0$, where it is very large. However, according to Equation (2.36), the function still possess a finite integral:

$$\int_{-\infty}^{\infty} \delta(x) dx = 1. \quad (2.38)$$

Thus, $\delta(x)$ has all of the required properties of a spike function. The one-dimensional spike function $\delta(x)$ is called the *Dirac delta function*, after Paul Dirac who invented it in 1927 while investigating quantum mechanics. The delta function is an example of what mathematicians call a *generalized function*; it is not well defined at $x = 0$, but its integral is nevertheless well defined. Consider the integral

$$\int_{-\infty}^{\infty} f(x) \delta(x) dx, \quad (2.39)$$

where $f(x)$ is a function that is well behaved in the vicinity of $x = 0$. Because the delta function is zero everywhere, apart from arbitrarily close to $x = 0$, it is clear that

$$\int_{-\infty}^{\infty} f(x) \delta(x) dx = f(0) \int_{-\infty}^{\infty} \delta(x) dx = f(0), \quad (2.40)$$

where use has been made of Equation (2.38). A simple change of variables allows us to define $\delta(x - x_0)$, which is a delta function centered on $x = x_0$. Equation (2.40) gives

$$\int_{-\infty}^{\infty} f(x) \delta(x - x_0) dx = f(x_0). \quad (2.41)$$

We actually require a three-dimensional delta function; that is, a function that is zero everywhere, apart from arbitrarily close to the origin, where it is very large, and whose volume integral

is unity. If we denote this function by $\delta(\mathbf{r})$ then it is easily seen that the three-dimensional delta function is the product of three one-dimensional delta functions:

$$\delta(\mathbf{r}) = \delta(x) \delta(y) \delta(z). \quad (2.42)$$

This function is clearly zero everywhere, except arbitrarily close the origin, where it is very large. But, is its volume integral unity? Let us integrate over a cube of dimension $2a$ that is centered on the origin, and aligned along the Cartesian axes. This volume integral is obviously separable, so that

$$\int \delta(\mathbf{r}) dV = \int_{-a}^a \delta(x) dx \int_{-a}^a \delta(y) dy \int_{-a}^a \delta(z) dz. \quad (2.43)$$

(See Section A.17.) The integral can be turned into an integral over all space by taking the limit $a \rightarrow \infty$. However, we know that, for one-dimensional delta functions, $\int_{-\infty}^{\infty} \delta(s) ds = 1$, so it follows from the previous equation that

$$\int \delta(\mathbf{r}) dV = 1, \quad (2.44)$$

which is the desired result. A simple generalization of previous arguments yields

$$\int f(\mathbf{r}) \delta(\mathbf{r}) dV = f(\mathbf{0}), \quad (2.45)$$

where $f(\mathbf{r})$ is any well-behaved scalar field. Finally, we can change variables and write

$$\delta(\mathbf{r} - \mathbf{r}') = \delta(x - x') \delta(y - y') \delta(z - z'), \quad (2.46)$$

which is a three-dimensional delta function centered on $\mathbf{r} = \mathbf{r}'$. It is easily demonstrated that

$$\int f(\mathbf{r}) \delta(\mathbf{r} - \mathbf{r}') dV = f(\mathbf{r}'). \quad (2.47)$$

Up to now, we have only considered volume integrals taken over all space. However, it should be obvious that the previous result also holds for integrals over any finite volume V that contains the point $\mathbf{r} = \mathbf{r}'$. Likewise, the integral is zero if V does not contain the point $\mathbf{r} = \mathbf{r}'$.

Let us now return to the problem in hand. The electric field generated by an electric charge q located at the origin has $\nabla \cdot \mathbf{E} = 0$ everywhere apart from the origin, and also satisfies

$$\int_V \nabla \cdot \mathbf{E} dV = \frac{q}{\epsilon_0} \quad (2.48)$$

for a spherical volume V centered on the origin. These two facts imply that

$$\nabla \cdot \mathbf{E} = \frac{q}{\epsilon_0} \delta(\mathbf{r}), \quad (2.49)$$

where use has been made of Equation (2.44).

Consider, again, an electric charge q located at the origin, and surrounded by a spherical surface S that is centered on the origin. We have seen that the flux of the electric field out of S is q/ϵ_0 .

Suppose that we now displace the surface S , so that it is no longer centered on the origin. What now is the flux of the electric field out of S ? We have

$$\oint_S \mathbf{E} \cdot d\mathbf{S} = \int_V \nabla \cdot \mathbf{E} dV \quad (2.50)$$

from the divergence theorem (see Section A.20), as well as Equation (2.49). From these two equations, it is clear that the flux of \mathbf{E} out of S is still q/ϵ_0 , as long as the displacement is not large enough that the origin is no longer enclosed by the sphere. Suppose that the surface S is not spherical, but is instead highly distorted. What now is the flux of \mathbf{E} out of S ? As before, the divergence theorem and Equation (2.49) tell us that the flux remains q/ϵ_0 , provided that the surface contains the origin. Moreover, this result is completely independent of the shape of S .

Let us try to extend the previous result. Consider N electric charges q_i located at displacements \mathbf{r}_i . A simple generalization of Equation (2.49) gives

$$\nabla \cdot \mathbf{E} = \sum_{i=1, N} \frac{q_i}{\epsilon_0} \delta(\mathbf{r} - \mathbf{r}_i). \quad (2.51)$$

Thus, Equation (2.50) and the previous equation imply that

$$\oint_S \mathbf{E} \cdot d\mathbf{S} = \int_V \nabla \cdot \mathbf{E} dV = \frac{Q}{\epsilon_0}, \quad (2.52)$$

where Q is the total charge enclosed by the surface S . This result is called *Gauss's law*, and does not depend on the shape of the surface. Note that the previous equation is analogous in form to the gravitational version of Gauss's law, (1.245). This is not surprising because, as we previously mentioned, Gauss's law holds for any inverse-square force law.

Suppose, finally, that instead of having a set of discrete electric charges, we have a continuous charge distribution described by a charge density $\rho(\mathbf{r})$. The charge contained in a small rectangular volume of dimensions dx , dy , and dz , located at displacement \mathbf{r} , is $Q = \rho(\mathbf{r}) dx dy dz$. However, if we integrate $\nabla \cdot \mathbf{E}$ over this volume element then we obtain

$$\nabla \cdot \mathbf{E} dx dy dz = \frac{Q}{\epsilon_0} = \frac{\rho dx dy dz}{\epsilon_0}, \quad (2.53)$$

where use has been made of Equation (2.52). Here, the volume element is assumed to be sufficiently small that $\nabla \cdot \mathbf{E}$ does not vary significantly across it. Thus, we get

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}. \quad (2.54)$$

Equation (2.54) is a differential equation that describes the electric field generated by a set of charges. We already know the solution to this equation when the charges are stationary; it is given by Equation (2.13),

$$\mathbf{E}(\mathbf{r}) = \frac{1}{4\pi \epsilon_0} \int_{V'} \rho(\mathbf{r}') \frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3} dV'. \quad (2.55)$$

Incidentally, Equations (2.54) and (2.55) can be reconciled provided

$$\nabla \cdot \left(\frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3} \right) = -\nabla^2 \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right) = 4\pi \delta(\mathbf{r} - \mathbf{r}'), \quad (2.56)$$

where use has been made of Equation (2.16). (See Section A.21.) It follows that

$$\begin{aligned} \nabla \cdot \mathbf{E}(\mathbf{r}) &= \frac{1}{4\pi \epsilon_0} \int \rho(\mathbf{r}') \nabla \cdot \left(\frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3} \right) dV' \\ &= \int \frac{\rho(\mathbf{r}')}{\epsilon_0} \delta(\mathbf{r} - \mathbf{r}') dV' = \frac{\rho(\mathbf{r})}{\epsilon_0}, \end{aligned} \quad (2.57)$$

which is the desired result. Here, use has been made of Equation (2.47).

Finally, the most general form of Gauss's law, Equation (2.52), is obtained by integrating Equation (2.54) over a volume V surrounded by a surface S , and making use of the divergence theorem:

$$\oint_S \mathbf{E} \cdot d\mathbf{S} = \frac{1}{\epsilon_0} \int_V \rho(\mathbf{r}) dV. \quad (2.58)$$

(See Section A.20.)

2.1.7 Applications of Gauss's Law

One particularly interesting application of Gauss's law is *Earnshaw's theorem*, which states that it is impossible for a collection of electrically charged particles to remain in static equilibrium solely under the influence of (classical) electrostatic forces. For instance, consider the motion of the i th particle in the electric field, $\mathbf{E}(\mathbf{r})$, generated by all of the other static particles. The equilibrium position of the i th particle corresponds to some point of displacement \mathbf{r}_i at which $\mathbf{E}(\mathbf{r}_i) = \mathbf{0}$, because this implies that the particle is not subject to an electrical force. By implication, \mathbf{r}_i does not correspond to the equilibrium displacement of any other particle in the system. However, in order for \mathbf{r}_i to be the displacement of a stable equilibrium point, the i th particle must experience a restoring force when its displacement deviates slightly from \mathbf{r}_i in any direction. Assuming that the i th particle is (say) positively charged, this implies that the electric field must be directed radially toward the point whose displacement is \mathbf{r}_i at all neighboring points. Hence, if we consider a small sphere centered on displacement \mathbf{r}_i then there must be a negative flux of \mathbf{E} through the surface of this sphere. According to Gauss's law, this necessitates the presence of a negative charge at displacement \mathbf{r}_i . However, there is no such charge at displacement \mathbf{r}_i . Hence, we conclude that \mathbf{E} cannot be directed radially toward the point whose displacement is \mathbf{r}_i at all neighboring points. In other words, there must be some neighboring points at which \mathbf{E} is directed away from the point whose displacement is \mathbf{r}_i . Hence, a positively charged particle placed at displacement \mathbf{r}_i can always escape by moving to such neighboring points. One corollary of Earnshaw's theorem is that classical electrostatics cannot account for the stability of atoms and molecules.

As an example of the use of Gauss's law, let us calculate the electric field generated by a spherically symmetric charge annulus of inner radius a , and outer radius b , centered on the origin,

and carrying a uniformly distributed electric charge Q . Now, by symmetry, we expect a spherically symmetric charge distribution to generate a spherically symmetric potential, $\phi(r)$, where r is a spherical polar coordinate. (See Section A.23.) It therefore follows from Equation (2.17) that the electric field is both spherically symmetric and radial; that is, $\mathbf{E} = E_r(r) \mathbf{e}_r$. Let us apply Gauss's law to an imaginary spherical surface, of radius r , centered on the origin. See Figure 2.4. Such a surface is generally known as a *Gaussian surface*. According to Gauss's law, (2.58), the flux of the electric field out of the surface is equal to the enclosed charge, divided by ϵ_0 . The flux is easy to calculate because the electric field is everywhere perpendicular to the surface. We obtain

$$4\pi r^2 E_r(r) = \frac{Q(r)}{\epsilon_0}, \quad (2.59)$$

where $Q(r)$ is the charge enclosed by a Gaussian surface of radius r . However, simple arguments involving proportion reveal that

$$Q(r) = \begin{cases} 0 & r < a \\ [(r^3 - a^3)/(b^3 - a^3)] Q & a \leq r \leq b \\ Q & b < r \end{cases} . \quad (2.60)$$

Hence,

$$E_r(r) = \begin{cases} 0 & r < a \\ [Q/(4\pi \epsilon_0 r^2)] [(r^3 - a^3)/(b^3 - a^3)] & a \leq r \leq b \\ Q/(4\pi \epsilon_0 r^2) & b < r \end{cases} . \quad (2.61)$$

The previous electric field distribution illustrates two important points. First, the electric field generated outside a spherically symmetric charge distribution is the same as that which would be generated if all of the charge in the distribution was concentrated at its center. Second, zero electric field is generated inside an empty cavity surrounded by a spherically symmetric charge distribution.

We can easily determine the electric potential associated with the electric field (2.61) using

$$\frac{d\phi(r)}{dr} = -E_r(r). \quad (2.62)$$

[See Equation (2.17).] The boundary conditions are that $\phi(\infty) = 0$, and that $\phi(r)$ is continuous at $r = a$ and $r = b$. (Of course, a discontinuous potential would lead to an infinite electric field, which is unphysical.) It follows that

$$\phi(r) = \begin{cases} [Q/(4\pi \epsilon_0)] (3/2) [(b^2 - a^2)/(b^3 - a^3)] & r < a \\ [Q/(4\pi \epsilon_0 r)] [(3b^3 r - r^3 - 2a^3)/2(b^3 - a^3)] & a \leq r \leq b \\ Q/(4\pi \epsilon_0 r) & b < r \end{cases} . \quad (2.63)$$

Hence, the work done in slowly moving a charge from infinity to the center of the distribution (which is minus the work done by the electric field) is

$$W = q [\phi(0) - \phi(\infty)] = \frac{qQ}{4\pi \epsilon_0} \frac{3}{2} \left(\frac{b^2 - a^2}{b^3 - a^3} \right). \quad (2.64)$$

[See Equation (2.23).]

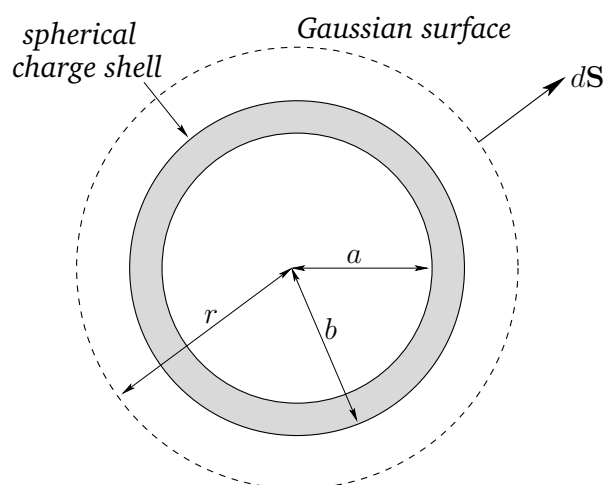


Figure 2.4: An example use of Gauss's law.

2.1.8 Electrostatic Energy

Consider a collection of N static point electric charges q_i located at displacements \mathbf{r}_i . What is the electrostatic energy stored in such a collection? In other words, how much work would we have to perform in order to assemble the charges, starting from an initial state in which they are all at rest and very widely separated?

We know that a static electric field is conservative, and can consequently be written in terms of a scalar potential:

$$\mathbf{E} = -\nabla\phi. \quad (2.65)$$

[See Equation (2.17).] We also know that the electrical force acting on a charge q located at displacement \mathbf{r} is written

$$\mathbf{f} = q\mathbf{E}(\mathbf{r}). \quad (2.66)$$

[See Equation (2.10).] The work that we would have to do against electrical forces in order to slowly move the charge from point P to point Q is simply

$$W = \int_P^Q (-\mathbf{f}) \cdot d\mathbf{r} = -q \int_P^Q \mathbf{E} \cdot d\mathbf{r} = q \int_P^Q \nabla\phi \cdot d\mathbf{r} = q [\phi(Q) - \phi(P)], \quad (2.67)$$

where $d\mathbf{r}$ is an element of the path taken between the two points. (See Section 1.3.2.) The negative sign in the previous expression comes about because we would have to exert a force $-\mathbf{f}$ on the charge, in order to counteract the force exerted by the electric field. Recall, finally, that the scalar potential field generated by a point charge q located at position \mathbf{r}' is

$$\phi(\mathbf{r}) = \frac{1}{4\pi\epsilon_0} \frac{q}{|\mathbf{r} - \mathbf{r}'|}. \quad (2.68)$$

[See Equation (2.21).]

Let us build up our collection of charges one by one. It takes no work to bring the first charge from infinity, because there is no electric field to fight against. Let us clamp this charge in position

at displacement \mathbf{r}_1 . In order to bring the second charge into position at displacement \mathbf{r}_2 , we have to do work against the electric field generated by the first charge. According to Equations (2.67) and (2.68), this work is given by

$$W_2 = \frac{1}{4\pi \epsilon_0} \frac{q_2 q_1}{|\mathbf{r}_2 - \mathbf{r}_1|}. \quad (2.69)$$

Let us now bring the third charge into position. Because electric fields and scalar potentials are superposable, the work done while moving the third charge from infinity to displacement \mathbf{r}_3 is simply the sum of the works done against the electric fields generated by charges 1 and 2, taken in isolation:

$$W_3 = \frac{1}{4\pi \epsilon_0} \left(\frac{q_3 q_1}{|\mathbf{r}_3 - \mathbf{r}_1|} + \frac{q_3 q_2}{|\mathbf{r}_3 - \mathbf{r}_2|} \right). \quad (2.70)$$

Thus, the total work done in assembling the collection of three charges is given by

$$W = \frac{1}{4\pi \epsilon_0} \left(\frac{q_2 q_1}{|\mathbf{r}_2 - \mathbf{r}_1|} + \frac{q_3 q_1}{|\mathbf{r}_3 - \mathbf{r}_1|} + \frac{q_3 q_2}{|\mathbf{r}_3 - \mathbf{r}_2|} \right). \quad (2.71)$$

This result can easily be generalized to a collection of N charges:

$$W = \frac{1}{4\pi \epsilon_0} \sum_{i=1, N} \sum_{j=1, N}^{j < i} \frac{q_i q_j}{|\mathbf{r}_i - \mathbf{r}_j|}. \quad (2.72)$$

The restriction that j must be less than i makes the previous summation rather messy. If we were to sum without restriction (other than $j \neq i$) then each pair of charges would be counted twice. It is convenient to do just this, and then to divide the result by two. Thus, we obtain

$$W = \frac{1}{2} \frac{1}{4\pi \epsilon_0} \sum_{i=1, N} \sum_{j=1, N}^{j \neq i} \frac{q_i q_j}{|\mathbf{r}_i - \mathbf{r}_j|}. \quad (2.73)$$

This is the electric potential energy (i.e., the difference between the total energy and the kinetic energy) of a collection of point electric charges. We can think of this quantity as the work required to bring stationary charges from infinity and assemble them in the required formation. Alternatively, it is the kinetic energy that would be released if the collection were dissolved, and the charges returned to infinity. But where is this potential energy stored? Let us investigate further.

Equation (2.73) can be written

$$W = \frac{1}{2} \sum_{i=1, N} q_i \phi_i, \quad (2.74)$$

where

$$\phi_i = \frac{1}{4\pi \epsilon_0} \sum_{j=1, N}^{j \neq i} \frac{q_j}{|\mathbf{r}_i - \mathbf{r}_j|} \quad (2.75)$$

is the scalar potential experienced by the i th charge due to the other charges in the distribution. [See Equation (2.20).]

Let us now consider the potential energy of a continuous charge distribution. It is tempting to write

$$W = \frac{1}{2} \int \rho \phi dV, \quad (2.76)$$

by analogy with Equations (2.74) and (2.75), where

$$\phi(\mathbf{r}) = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} dV' \quad (2.77)$$

is the familiar scalar potential generated by a continuous charge distribution of charge density $\rho(\mathbf{r})$ [see Equation (2.18)], and where the volume integrals are over all space. Let us try this scheme out. We know from Equation (2.54) that

$$\rho = \epsilon_0 \nabla \cdot \mathbf{E}, \quad (2.78)$$

so Equation (2.76) can be written

$$W = \frac{\epsilon_0}{2} \int \phi \nabla \cdot \mathbf{E} dV. \quad (2.79)$$

Now,

$$\nabla \cdot (\mathbf{E} \phi) \equiv \phi \nabla \cdot \mathbf{E} + \mathbf{E} \cdot \nabla \phi. \quad (2.80)$$

(See Section A.24.) However, $\nabla \phi = -\mathbf{E}$, so we obtain

$$W = \frac{\epsilon_0}{2} \left[\int \nabla \cdot (\mathbf{E} \phi) dV + \int E^2 dV \right] \quad (2.81)$$

Application of the divergence theorem (see Section A.20) gives

$$W = \frac{\epsilon_0}{2} \left(\oint_S \phi \mathbf{E} \cdot d\mathbf{S} + \int_V E^2 dV \right), \quad (2.82)$$

where V is some volume that encloses all of the charges, and S is its bounding surface. Let us assume that V is a sphere, centered on the origin, and let us take the limit in which the radius r of this sphere goes to infinity. We know that, in general, the electric field at large distances from a bounded charge distribution looks like the field of a point charge, and, therefore, falls off like $1/r^2$. Likewise, the potential falls off like $1/r$. (See Section 2.1.7.) However, the surface area of the sphere increases like r^2 . Hence, it is clear that, in the limit as $r \rightarrow \infty$, the surface integral in Equation (2.82) falls off like $1/r$, and is consequently zero. Thus, Equation (2.82) reduces to

$$W = \frac{\epsilon_0}{2} \int E^2 dV, \quad (2.83)$$

where the volume integral is over all space. This is a very interesting result. It tells us that the potential energy of a continuous charge distribution is stored in the electric field generated by that distribution. Of course, we now have to assume that an electric field possesses an energy density

$$U = \frac{\epsilon_0}{2} E^2. \quad (2.84)$$

Incidentally, the fact that an electric field possess an energy density demonstrates that it has a real physical existence, and is not just an aid to the calculation of electrostatic forces.

We can easily check that Equation (2.83) is correct. Suppose that we have an electric charge Q that is uniformly distributed within a sphere of radius a centered on the origin. Let us imagine building up this charge distribution from a succession of thin spherical layers of infinitesimal thickness. At each stage, we gather a small amount of charge dq from infinity, and spread it over the surface of the sphere in a thin layer extending from r to $r + dr$. We continue this process until the final radius of the sphere is a . If $q(r)$ is the sphere's charge when it has attained radius r then the work done in bringing a charge dq to its surface is

$$dW = \frac{1}{4\pi\epsilon_0} \frac{q(r) dq}{r}. \quad (2.85)$$

This follows from Equation (2.69), because the electric field generated outside a spherical charge distribution is the same as that of a point charge $q(r)$ located at its geometric center ($r = 0$). (See Section 2.1.7.) If the constant charge density of the sphere is ρ then

$$q(r) = \frac{4\pi}{3} r^3 \rho, \quad (2.86)$$

and

$$dq = 4\pi r^2 \rho dr. \quad (2.87)$$

Thus, Equation (2.85) becomes

$$dW = \frac{4\pi}{3\epsilon_0} \rho^2 r^4 dr. \quad (2.88)$$

The total work needed to build up the sphere from zero radius to radius a is plainly

$$W = \frac{4\pi}{3\epsilon_0} \rho^2 \int_0^a r^4 dr = \frac{4\pi}{15\epsilon_0} \rho^2 a^5. \quad (2.89)$$

This can also be written in terms of the total charge $Q = (4\pi/3) a^3 \rho$ as

$$W = \frac{3}{5} \frac{Q^2}{4\pi\epsilon_0 a}. \quad (2.90)$$

Now that we have evaluated the potential energy of a spherical charge distribution by the direct method, let us work it out using Equation (2.83). We shall assume that the electric field is both radial and spherically symmetric, so that $\mathbf{E} = E_r(r) \mathbf{e}_r$. Here, r is a standard spherical polar coordinate. (See Section A.23.) Application of Gauss's law,

$$\oint_S \mathbf{E} \cdot d\mathbf{S} = \frac{1}{\epsilon_0} \int_V \rho dV, \quad (2.91)$$

where V is a sphere of radius r , centered on the origin, gives

$$E_r(r) = \frac{Q}{4\pi\epsilon_0 a^3} \frac{r}{a} \quad (2.92)$$

for $r < a$, and

$$E_r(r) = \frac{Q}{4\pi \epsilon_0 r^2} \quad (2.93)$$

for $r \geq a$. Equations (2.83), (2.92), and (2.93) yield

$$W = \frac{Q^2}{8\pi \epsilon_0} \left(\frac{1}{a^6} \int_0^a r^4 dr + \int_a^\infty \frac{dr}{r^2} \right), \quad (2.94)$$

which reduces to

$$W = \frac{Q^2}{8\pi \epsilon_0 a} \left(\frac{1}{5} + 1 \right) = \frac{3}{5} \frac{Q^2}{4\pi \epsilon_0 a}. \quad (2.95)$$

Thus, Equation (2.83) gives the correct answer.

The reason that we have checked Equation (2.83) so carefully is that, on close inspection, it is found to be inconsistent with Equation (2.74), from which it was supposedly derived. For instance, the energy given by Equation (2.83) is manifestly positive definite, whereas the energy given by Equation (2.74) can be negative (it is certainly negative for a collection of two point charges of opposite sign). The inconsistency was introduced into our analysis when we replaced Equation (2.75) by Equation (2.77). In Equation (2.75), the self-interaction of the i th charge with its own electric field is specifically excluded, whereas it is included in Equation (2.77). Thus, the potential energies (2.74) and (2.83) are different because in the former we start from ready-made point charges, whereas in the latter we build up the whole charge distribution from scratch. Hence, if we were to work out the potential energy of a point charge distribution using Equation (2.83) then we would obtain the energy (2.74) plus the energy required to assemble the point charges. What is the energy required to assemble a point electric charge? In fact, it is infinite. To see this, let us suppose, for the sake of argument, that our point charges actually consist of electric charge uniformly distributed in small spheres of radius b . According to Equation (2.90), the energy required to assemble the i th point charge is

$$W_i = \frac{3}{5} \frac{q_i^2}{4\pi \epsilon_0 b}. \quad (2.96)$$

We can think of this as the self-energy of the i th charge. Thus, we can write

$$W = \frac{\epsilon_0}{2} \int E^2 dV = \frac{1}{2} \sum_{i=1,N} q_i \phi_i + \sum_{i=1,N} W_i \quad (2.97)$$

which enables us to reconcile Equations (2.74) and (2.83). Unfortunately, if our point charges really are point charges then $b \rightarrow 0$, and the self-energy of each charge becomes infinite. Thus, the potential energies predicted by Equations (2.74) and (2.83) differ by an infinite amount. What does this all mean? We have to conclude that the idea of locating electrostatic potential energy in the electric field runs into conceptual difficulties in the presence of point electric charges. One way out of this dilemma would be to say that elementary electric charges, such as protons and electrons, are not point objects, but instead have finite spatial extents. Regrettably, although protons have finite spatial extents (of about 10^{-15} m), electrons really do seem to be point objects.

2.1.9 Poisson's Equation

We have seen that the electric field generated by a set of stationary charges can be written as the gradient of a scalar potential, so that

$$\mathbf{E} = -\nabla\phi. \quad (2.98)$$

[See Equation (2.17).] The previous equation can be combined with the field equation (2.54) to give a partial differential equation for the scalar potential:

$$\nabla^2\phi = -\frac{\rho}{\epsilon_0}. \quad (2.99)$$

(See Section A.21.) This equation is known as *Poisson's equation*. (See Section 1.8.5.)

2.1.10 Uniqueness Theorem

Consider a volume V bounded by some surface S . Suppose that we are given the electric charge density ρ throughout V , and the (not necessarily constant) value of the scalar potential ϕ_S on S . Is this sufficient information for Poisson's equation to uniquely specify the scalar potential throughout V ? Suppose, for the sake of argument, that the solution is not unique. Let there be two different potentials ϕ_1 and ϕ_2 that both satisfy

$$\nabla^2\phi_1 = -\frac{\rho}{\epsilon_0}, \quad (2.100)$$

$$\nabla^2\phi_2 = -\frac{\rho}{\epsilon_0} \quad (2.101)$$

throughout V , and

$$\phi_1 = \phi_S, \quad (2.102)$$

$$\phi_2 = \phi_S \quad (2.103)$$

on S . We can form the difference between these two potentials:

$$\phi_3 = \phi_1 - \phi_2. \quad (2.104)$$

The potential ϕ_3 clearly satisfies

$$\nabla^2\phi_3 = 0 \quad (2.105)$$

throughout V , and

$$\phi_3 = 0 \quad (2.106)$$

on S .

Now,

$$\nabla \cdot (\phi_3 \nabla\phi_3) \equiv (\nabla\phi_3)^2 + \phi_3 \nabla^2\phi_3. \quad (2.107)$$

(See Section A.24.) Thus, making use of the divergence theorem,

$$\int_V [(\nabla\phi_3)^2 + \phi_3 \nabla^2\phi_3] dV = \oint_S \phi_3 \nabla\phi_3 \cdot d\mathbf{S}. \quad (2.108)$$

(See Section A.20.) But, $\nabla^2\phi_3 = 0$ throughout V , and $\phi_3 = 0$ on S , so the previous equation reduces to

$$\int_V (\nabla\phi_3)^2 dV = 0. \quad (2.109)$$

Note that $(\nabla\phi_3)^2$ is a positive definite quantity. The only way in which the volume integral of a positive definite quantity can be zero is if that quantity itself is zero throughout the volume. This is not necessarily the case for a non-positive definite quantity, because we could have positive and negative contributions from various regions inside the volume that cancel one another out. Thus, because $(\nabla\phi_3)^2 \equiv \nabla\phi_3 \cdot \nabla\phi_3$ is positive definite, it is zero throughout V . It follows that $\nabla\phi_3 = \mathbf{0}$ throughout V , and, hence, that

$$\phi_3 = \text{constant} \quad (2.110)$$

throughout V . However, we know that $\phi_3 = 0$ on S , so we get

$$\phi_3 = 0 \quad (2.111)$$

throughout V . In other words,

$$\phi_1 = \phi_2 \quad (2.112)$$

throughout V and on S . Our initial supposition that ϕ_1 and ϕ_2 are two different solutions of Poisson's equation, satisfying the same boundary conditions, turns out to be incorrect. Hence, we deduce that the solutions to Poisson's equation in a volume bounded by a surface on which the electric potential is specified are unique. This important result is known as the *uniqueness theorem*.

2.1.11 Ohm's Law

A *conductor* is a medium that contains free electric charges (usually electrons) that acquire a net drift velocity in the presence of an applied electric field, giving rise to an *electric current* flowing in the same direction as the field. The well-known relationship between the current and the voltage in a typical conductor is given by *Ohm's law*:

$$V = IR, \quad (2.113)$$

where V is the voltage drop across a conductor of electrical resistance R through which a current I flows. Incidentally, the unit of electric current is the *ampere* (or amp) (A), which is equivalent to a coulomb per second. Furthermore, the unit of electrical resistance is the *ohm* (Ω), which is equivalent to a volt per ampere.

Let us generalize Ohm's law so that it is expressed in terms of the electric field, \mathbf{E} , and *current density*, \mathbf{j} , at a given point inside the conductor, rather than the global quantities V and I . Here, the magnitude of the current density vector, \mathbf{j} , measures the amount of current flowing per unit time per unit cross-sectional area, whereas the direction of the vector specifies the direction of the current flow. Consider a length l of a conductor of uniform cross-sectional area A through which a net electric current I flows. In general, we expect the electrical resistance of the conductor to be proportional to its length, l , and inversely proportional to its cross-sectional area, A (i.e., we expect

it to be harder to push an electrical current down a long rather than a short wire, and easier to push an electrical current down a wide rather than a narrow conducting channel.) Thus, we can write

$$R = \eta \frac{l}{A}. \quad (2.114)$$

Here, the constant η is called the *resistivity* of the conducting medium, and is measured in units of ohm-meters. Hence, Ohm's law becomes

$$V = \eta \frac{l}{A} I. \quad (2.115)$$

However, $I/A = j_z$ (supposing that the conductor is aligned along the z -axis) and $V/l = E_z$ [see Equation (2.17)], so the previous equation reduces to

$$E_z = \eta j_z. \quad (2.116)$$

Because there is nothing special about the z -axis (in an isotropic conducting medium), the previous formula immediately generalizes to

$$\mathbf{E} = \eta \mathbf{j}. \quad (2.117)$$

This is the most fundamental form of Ohm's law.

It is fairly easy to account for the previous equation at the microscopic level. Consider a metal that has n_e free electrons per unit volume. Of course, the metal also has a fixed lattice of metal ions whose charge per unit volume is equal and opposite to that of the free electrons, rendering the medium electrically neutral. In the presence of an electric field \mathbf{E} , a given free electron is subject to an electrical force $\mathbf{f} = -e \mathbf{E}$ [see Equation (2.10)], and therefore accelerates (from rest at $t = 0$) such that its drift velocity is written $\mathbf{v} = -(e/m_e) t \mathbf{E}$, where $-e$ is the electron charge, and m_e the electron mass. Suppose that, on average, a drifting electron collides with a metal ion once every τ seconds. Given that a metal ion is much more massive than an electron, we expect a free electron to lose all of the momentum it had previously acquired from the electric field during such a collision. It follows that the mean drift velocity of the free electrons is $\langle \mathbf{v} \rangle = -(e \tau / 2 m_e) \mathbf{E}$. Hence, the mean current density is

$$\mathbf{j} = -n_e e \langle \mathbf{v} \rangle = \frac{n_e e^2 \tau}{2 m_e} \mathbf{E}. \quad (2.118)$$

Thus, we can account for Equation (2.117), as long as the resistivity takes the form

$$\eta = \frac{2 m_e}{n_e e^2 \tau}. \quad (2.119)$$

We conclude that the resistivity of a typical conducting medium is determined by the number density of free electrons, as well as the mean collision rate of these electrons with the fixed ions.

A free charge q that moves through a voltage drop V acquires an energy qV from the electric field. (See Section 2.1.5.) In a conducting medium, this energy is dissipated as heat (the conversion to heat takes place each time a free charge collides with a fixed ion). This particular type of heating is called *ohmic heating* or *Joule heating*. Suppose that N charges per unit time pass through a

conductor. The current flowing through the conductor is obviously $I = Nq$. The total energy gained by the charges, which appears as heat inside the conductor, is

$$P = NqV = IV \quad (2.120)$$

per unit time. Thus, the heating power is

$$P = IV = I^2 R = \frac{V^2}{R}. \quad (2.121)$$

Equations (2.120) and (2.121) generalize to

$$P = \mathbf{j} \cdot \mathbf{E} = \eta j^2, \quad (2.122)$$

where P is now the power dissipated per unit volume inside the conducting medium.

2.1.12 Ideal Conductors

Most electrical conductors obey Ohm's law, and are termed *ohmic* conductors. Suppose that we apply an electric field to an ohmic conductor of very low resistivity. What is going to happen? According to Equation (2.117), the electric field drives very large currents inside the conductor. These currents will redistribute the electric charge within the conductor until the original electric field is canceled out. At this point, the currents stop flowing. It might be objected that the currents could keep flowing in closed loops. According to Ohm's law, this would require a non-zero *electromotive force* (emf), $\oint \mathbf{E} \cdot d\mathbf{r}$, acting around each loop (unless the conductor is a superconductor, with $\eta = 0$). However, we know that in a steady state

$$\oint_C \mathbf{E} \cdot d\mathbf{r} = 0 \quad (2.123)$$

around any closed loop C . [See Equation (2.24).] This proves that a steady-state emf acting around a closed loop inside a conductor is impossible. The only other alternative is

$$\mathbf{j} = \mathbf{E} = \mathbf{0} \quad (2.124)$$

everywhere inside the conductor. It immediately follows from the field equation $\nabla \cdot \mathbf{E} = \rho/\epsilon_0$ [see Equation (2.54)] that

$$\rho = 0. \quad (2.125)$$

We conclude that there is zero net electric charge in the interior of an ideal conductor. But, how can a conductor cancel out an applied electric field if it contains no internal electric charge? The answer is that the requisite charges reside on the surface of the conductor. (In reality, the charges lie within one or two atomic layers of the surface.)

Now, the difference in scalar potential between two points P and Q is simply

$$\phi(Q) - \phi(P) = \int_P^Q \nabla\phi \cdot d\mathbf{r} = - \int_P^Q \mathbf{E} \cdot d\mathbf{r}. \quad (2.126)$$

[See Equation (2.17) and Section A.18.] However, if points P and Q both lie inside the same conductor then it is clear from Equations (2.124) and (2.126) that the potential difference between P and Q is zero. This is true no matter where P and Q are situated inside the conductor, so we conclude that the scalar potential must be uniform inside a conductor. One corollary of this fact is that the surface of a conductor is an *equipotential* (i.e., $\phi = \text{constant}$) surface.

We have demonstrated that the electric field inside a conductor is zero. We can also demonstrate that the field within an empty cavity lying inside a conductor is zero, provided that there are no charges within the cavity. Let V be the cavity in question, and let S be its bounding surface. Because there are no electric charges within the cavity, the electric potential, ϕ , inside the cavity satisfies

$$\nabla^2\phi = 0. \quad (2.127)$$

[See Equation (2.99).] However, because S corresponds to the inner surface of the conductor that surrounds the cavity, S is an equipotential surface. In other words, the electric potential on S takes a constant value, ϕ_S (say). So, we need to solve a simplified version of Poisson's equation, (2.127), throughout V , subject to the boundary condition that $\phi = \phi_S$ on S . One obvious solution to this problem is $\phi = \phi_S$ throughout V and on S . However, we showed in Section 2.1.10 that the solutions to Poisson's equation in a volume surrounded by a surface on which the potential is specified are unique. Thus, $\phi = \phi_S$ throughout V and on S is the only solution to the problem. It follows that the electric field $\mathbf{E} = -\nabla\phi$ is zero throughout the cavity. [See Equation (2.17).]

We have shown that if a charge-free cavity is completely enclosed by a conductor then no stationary distribution of charges outside the conductor can ever produce any electric fields inside the cavity. It follows that we can shield a sensitive piece of electrical equipment from stray external electric fields by placing it inside a metal can. In fact, a wire mesh cage will do, as long as the mesh spacing is not too wide. Such a cage is known as a *Faraday cage*.

Consider a small region lying on the surface of a conductor. Suppose that the local surface electric charge density is σ , and that the electric field just outside the conductor is \mathbf{E} . Note that this field must be directed normal to the surface of the conductor. Any parallel component would be shorted out by surface currents. Another way of saying this is that the surface of a conductor is an equipotential. We know that $\nabla\phi$ is always perpendicular to an equipotential (see Section A.18), so $\mathbf{E} = -\nabla\phi$ [see Equation (2.17)] must be locally perpendicular to a conducting surface. Let us use Gauss's law [see Equation (2.58)],

$$\oint_S \mathbf{E} \cdot d\mathbf{S} = \frac{1}{\epsilon_0} \int_V \rho dV, \quad (2.128)$$

where the volume V is a so-called *Gaussian pill-box*. See Figure 2.5. A Gaussian pill-box is a volume of space whose shape is similar to an old-fashioned pill-box (or a modern pizza box). Let the two flat ends of the pill-box be aligned parallel to the surface of the conductor, with the surface running between them, and let the comparatively short sides be perpendicular to the surface. It is clear that \mathbf{E} is parallel to the sides of the box, so the sides make no contribution to the surface integral. The end of the box that lies inside the conductor also makes no contribution, because $\mathbf{E} = \mathbf{0}$ inside a conductor. Thus, the only non-zero contribution to the surface integral comes from the end lying in free space. This contribution is simply $E_{\perp} A$, where E_{\perp} denotes an outward

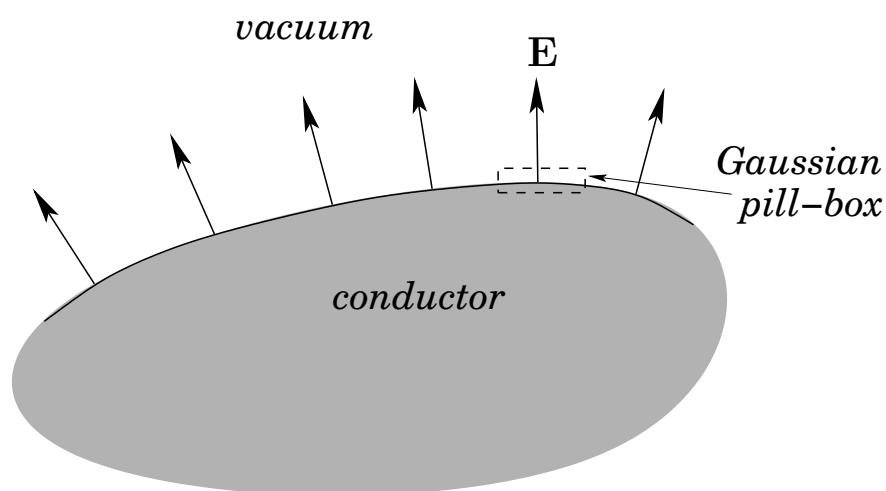


Figure 2.5: The surface of a conductor.

pointing (from the conductor) normal electric field, and A is the cross-sectional area of the box. The charge enclosed by the box is simply σA , from the definition of a surface charge density. Thus, Gauss's law yields

$$E_{\perp} = \frac{\sigma}{\epsilon_0} \quad (2.129)$$

as the relationship between the normal electric field immediately outside a conductor and the surface charge density.

Let us look at the electric field generated by a sheet charge distribution a little more carefully. Suppose that the charge per unit area is σ . By symmetry, we expect the field generated below the sheet to be the mirror image of that above the sheet (at least, locally). Thus, if apply Gauss's law to a pill-box of cross-sectional area A , as shown in Figure 2.6, then the two ends both contribute $E_{\text{sheet}} A$ to the surface integral, where E_{sheet} is the normal electric field generated above and below the sheet. The charge enclosed by the pill-box is just σA . Thus, Gauss's law yields a symmetric electric field

$$E_{\text{sheet}} = \begin{cases} +\sigma/(2\epsilon_0) & \text{above} \\ -\sigma/(2\epsilon_0) & \text{below} \end{cases} \quad (2.130)$$

So, how do we get the asymmetric electric field of a conducting surface, which is zero immediately below the surface (i.e., inside the conductor) and non-zero immediately above it? Clearly, we have to add in an external field (i.e., a field that is not generated locally by the sheet charge). The requisite field is

$$E_{\text{ext}} = \frac{\sigma}{2\epsilon_0} \quad (2.131)$$

both above and below the charge sheet. The total field is the sum of the field generated locally by the charge sheet and the external field. Thus, we obtain

$$E_{\text{total}} = \begin{cases} +\sigma/\epsilon_0 & \text{above} \\ 0 & \text{below} \end{cases} \quad (2.132)$$

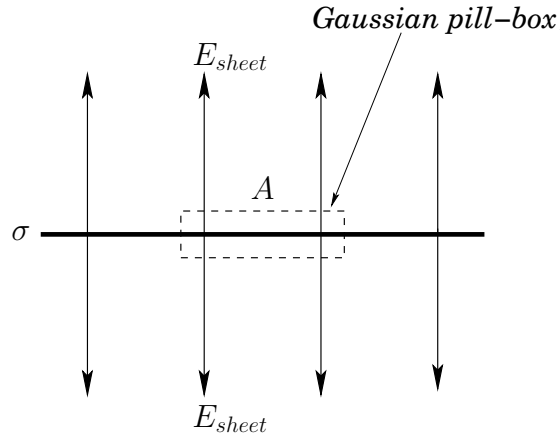


Figure 2.6: The electric field of a sheet charge.

which is in agreement with Equation (2.129). Now, the external electric field exerts a force on the charge sheet. Of course, the field generated locally by the sheet itself cannot exert a local force (i.e., by Newton's third law of motion, the charge sheet cannot locally exert a force on itself; see Section 1.2.4). Thus, the force per unit area acting on the surface of a conductor always acts outward, and is given by

$$p = \sigma E_{\text{ext}} = \frac{\sigma^2}{2 \epsilon_0}. \quad (2.133)$$

[See Equation (2.10).] We conclude that there is an *electrostatic pressure* acting on any charged conductor. This effect can be observed by charging up soap bubbles; the additional electrostatic pressure eventually causes them to burst.

Making use of Equations (2.131) and (2.133), the electrostatic pressure acting at the surface of a conductor can also be written

$$p = \frac{\epsilon_0}{2} E_{\perp}^2, \quad (2.134)$$

where E_{\perp} is the electric field-strength immediately above the surface of the conductor. Note that, according to the previous formula, the electrostatic pressure is equivalent to the energy density of the electric field immediately outside the conductor. [See Equation (2.84).] This is not a coincidence. Suppose that the conductor expands normally by an average distance dx , due to the electrostatic pressure. The electric field is excluded from the region into which the conductor expands. The volume of this region is $dV = A dx$, where A is the surface area of the conductor. Thus, the energy of the electric field decreases by an amount $dE = U dV = (\epsilon_0/2) E_{\perp}^2 dV$, where U is the energy density of the field. This decrease in energy can be ascribed to the work that the field does on the conductor in order to make it expand. This work is $dW = p A dx$, where p is the force per unit area that the field exerts on the conductor. Thus, $dE = dW$, from energy conservation, giving

$$p = \frac{\epsilon_0}{2} E_{\perp}^2. \quad (2.135)$$

Incidentally, this technique for calculating a force, given an expression for the energy of a system as a function of some adjustable parameter, is called *the principle of virtual work*.

2.1.13 Capacitors

It is clear that we can store electric charge on the surface of a conductor. However, electric fields will be generated immediately above this surface. Now, the conductor can only successfully store charge if it is electrically insulated from its surroundings. Of course, air is a very good electrical insulator. Unfortunately, air ceases to be an insulator when the electric field-strength through it exceeds some critical value which is about $E_{\text{crit}} \sim 10^6$ volts per meter. This phenomenon, which is known as *breakdown*, is associated with the formation of sparks. The most well-known example of the breakdown of air is during a lightning strike. Thus, a good charge-storing device is one that holds a relatively large amount of charge, but only generates relatively small external electric fields (so as to avoid breakdown). Such a device is called a *capacitor*.

Consider two thin, parallel, conducting plates of cross-sectional area A that are separated by a small distance d (i.e., $d \ll \sqrt{A}$). Suppose that each plate carries an equal and opposite charge $\pm Q$ (where $Q > 0$). We expect this charge to spread evenly over the plates to give an effective sheet charge density $\pm\sigma = Q/A$ on each plate. Suppose that the upper plate carries a positive charge and that the lower carries a negative charge. According to Equation (2.130), the field generated by the upper plate is normal to the plate and of magnitude

$$E_{\text{upper}} = \begin{cases} +\sigma/(2\epsilon_0) & \text{above} \\ -\sigma/(2\epsilon_0) & \text{below} \end{cases} \quad (2.136)$$

Likewise, the field generated by the lower plate is

$$E_{\text{lower}} = \begin{cases} -\sigma/(2\epsilon_0) & \text{above} \\ +\sigma/(2\epsilon_0) & \text{below} \end{cases} \quad (2.137)$$

Note that we are neglecting any leakage of the field at the edges of the plates. This is reasonable provided that the plates are relatively closely spaced. The total electric field is the sum of the two fields generated by the upper and lower plates. Thus, the net field is normal to the plates, and of magnitude

$$E_{\perp} = \begin{cases} \sigma/\epsilon_0 & \text{between} \\ 0 & \text{otherwise} \end{cases} \quad (2.138)$$

See Figure 2.7. Because the electric field between the plates is uniform, the potential difference between the plates is simply

$$V = E_{\perp} d = \frac{\sigma d}{\epsilon_0}. \quad (2.139)$$

[See Equation (2.17).]

It is conventional to measure the capacity of a conductor, or set of conductors, to store electric charge, but generate small external electric fields, in terms of a parameter called *capacitance*. This parameter is usually denoted C . The capacitance of a charge storing device is simply the ratio of the charge stored to the potential difference generated by this charge:

$$C = \frac{Q}{V}. \quad (2.140)$$

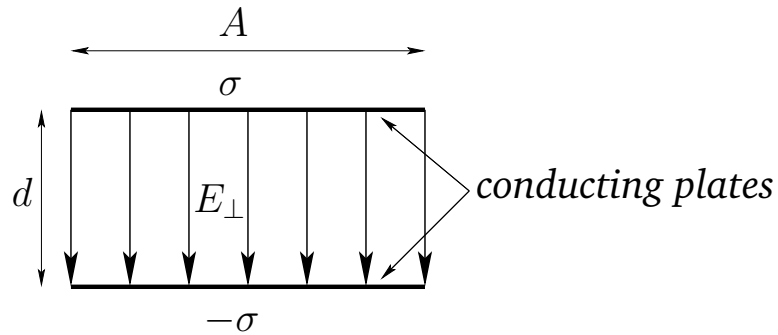


Figure 2.7: The electric field of a parallel plate capacitor.

Clearly, a good charge storing device has a high capacitance. Incidentally, capacitance is measured in *farads* (F), which are equivalent to coulombs per volt. This is a rather unwieldy unit, because capacitors in electrical circuits typically have capacitances that are only about one millionth of a farad.

For a parallel plate capacitor, we have

$$C = \frac{\sigma A}{V} = \frac{\epsilon_0 A}{d}. \quad (2.141)$$

Note that the capacitance only depends on geometric quantities, such as the area and spacing of the plates. This is a consequence of the superposability of electric fields. If we double the charge on a set of conductors then we double the electric fields generated around them, and we, therefore, double the potential difference between the conductors. Thus, the potential difference between the conductors is always directly proportional to the charge on the conductors. Moreover, the constant of proportionality (the inverse of the capacitance) can only depend on geometry.

Suppose that the charge $\pm Q$ on each plate of a parallel plate capacitor is built up gradually by transferring small amounts of charge from one plate to another. If the instantaneous charge on the plates is $\pm q$, and an infinitesimal amount of positive charge dq is transferred from the negatively charged to the positively charge plate, then the work done is $dW = V dq = q dq/C$, where V is the instantaneous voltage difference between the plates. (See Section 2.1.5.) Note that the voltage difference is such that it opposes any increase in the charge on either plate. The total work done in charging the capacitor is

$$W = \frac{1}{C} \int_0^Q q dq = \frac{Q^2}{2C} = \frac{1}{2} C V^2, \quad (2.142)$$

where use has been made of Equation (2.140). The energy stored in the capacitor is the same as the work required to charge up the capacitor. Thus, the stored energy is

$$W = \frac{1}{2} C V^2. \quad (2.143)$$

This is a general result that holds for all types of capacitor.

The energy of a charged parallel plate capacitor is actually stored in the electric field generated between the plates. This field is of approximately constant magnitude $E_{\perp} = V/d$, and occupies

a region of volume $A d$. Thus, given the energy density of an electric field, $U = (\epsilon_0/2) E^2$ [see Equation (2.84)], the energy stored in the electric field is

$$W = \frac{\epsilon_0}{2} \frac{V^2}{d^2} A d = \frac{1}{2} C V^2, \quad (2.144)$$

where use has been made of Equation (2.141). Note that Equations (2.142) and (2.144) agree with one another. The fact that the energy of a capacitor is stored in its electric field is also a general result.

The idea, that we discussed in the previous section, that an electric field exerts a negative pressure $(\epsilon_0/2) E_{\perp}^2$ on conductors immediately suggests that the two plates in a parallel plate capacitor attract one another with a mutual force

$$F = \frac{\epsilon_0}{2} E_{\perp}^2 A = \frac{1}{2} \frac{C V^2}{d}. \quad (2.145)$$

It is not actually necessary to have two oppositely charged conductors in order to make a capacitor. Consider an isolated conducting sphere of radius a , centered on the origin, that carries an electric charge Q . The spherically symmetric, radial electric field generated outside the sphere is given by

$$E_r(r > a) = \frac{Q}{4\pi \epsilon_0 r^2}, \quad (2.146)$$

and the associated electric potential is

$$\phi(r > a) = \frac{Q}{4\pi \epsilon_0 r}. \quad (2.147)$$

(See Section 2.1.7.) Here, r is a spherical polar coordinate. (See Section A.23.) It follows that the potential difference between the sphere and infinity—or, more realistically, some large, relatively distant reservoir of charge such as the Earth—is

$$V = \frac{Q}{4\pi \epsilon_0 a}. \quad (2.148)$$

Thus, the capacitance of the sphere is

$$C = \frac{Q}{V} = 4\pi \epsilon_0 a. \quad (2.149)$$

The energy of a spherical capacitor when it carries a charge Q is again given by $(1/2) C V^2$. It can easily be demonstrated that this is equivalent to the energy contained in the electric field surrounding the capacitor.

Suppose that we have two spheres of radii a and b , respectively, that are connected by a long electric wire. See Figure 2.8. The wire allows electric charge to move back and forth between the spheres until they reach the same potential (with respect to infinity). Let Q_a be the charge on the first sphere, and Q_b the charge on the second sphere. Of course, the total charge $Q_0 = Q_a + Q_b$

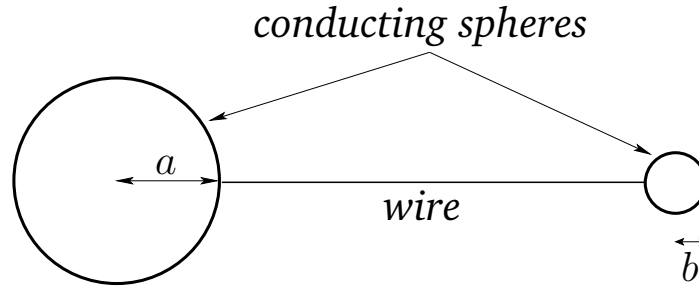


Figure 2.8: Two conducting spheres connected by a wire.

carried by the two spheres is a conserved quantity. It follows from Equation (2.148) that if the spheres are at the same potential then

$$\frac{Q_a}{Q_0} = \frac{a}{a+b}, \quad (2.150)$$

$$\frac{Q_b}{Q_0} = \frac{b}{a+b}. \quad (2.151)$$

Note that if one sphere is much smaller than the other one—for instance, if $b \ll a$ —then the large sphere grabs most of the charge; that is,

$$\frac{Q_a}{Q_b} \simeq \frac{a}{b} \gg 1. \quad (2.152)$$

The ratio of the electric fields generated just above the surfaces of the two spheres follows from Equations (2.146) and (2.152):

$$\frac{E_b}{E_a} \simeq \frac{a}{b}. \quad (2.153)$$

Note that if $b \ll a$ then the field just above the smaller sphere is far larger than that above the larger sphere. Equation (2.153) is a simple example of a far more general rule; namely, the electric field directly above some point on the surface of a conductor is inversely proportional to the local radius of curvature of the surface.

It is clear that if we wish to store significant amounts of charge on a conductor then the surface of the conductor must be made as smooth as possible. Any sharp spikes on the surface will inevitably have comparatively small radii of curvature. Intense local electric fields are thus generated around such spikes. These fields can easily exceed the critical field for the breakdown of air, leading to sparking and the eventual loss of the charge on the conductor. Sparking can also be very destructive, because the associated electric currents flow through very localized regions, giving rise to intense ohmic heating.

As a final example, consider two co-axial conducting cylinders of radii a and b , where $a < b$. Suppose that the charge per unit length carried by the outer and inner cylinders is $+\lambda$ and $-\lambda$, respectively. We can safely assume that $\mathbf{E} = E_r(r) \mathbf{e}_r$, by symmetry (adopting standard cylindrical

polar coordinates). (See Section A.23.) Let us apply Gauss's law (see Section 2.4) to a cylindrical surface of radius r , co-axial with the conductors, and of length l . For $a < r < b$, we find that

$$2\pi r l E_r(r) = \frac{\lambda l}{\epsilon_0}, \quad (2.154)$$

so that

$$E_r = \frac{\lambda}{2\pi \epsilon_0 r} \quad (2.155)$$

for $a < r < b$. It is fairly obvious that $E_r = 0$ if r is not in the range a to b . The potential difference between the inner and outer cylinders is [see Equation (2.17)]

$$V = - \int_{\text{outer}}^{\text{inner}} \mathbf{E} \cdot d\mathbf{r} = \int_{\text{inner}}^{\text{outer}} \mathbf{E} \cdot d\mathbf{r} = \int_a^b E_r dr = \frac{\lambda}{2\pi \epsilon_0} \int_a^b \frac{dr}{r}, \quad (2.156)$$

so

$$V = \frac{\lambda}{2\pi \epsilon_0} \ln \left(\frac{b}{a} \right). \quad (2.157)$$

Thus, the capacitance per unit length of the two cylinders is

$$C = \frac{\lambda}{V} = \frac{2\pi \epsilon_0}{\ln(b/a)}. \quad (2.158)$$

2.1.14 Method of Images

Suppose that a point electric charge q is located a distance d from an infinite, grounded (i.e., held at zero potential), conducting plate. See Figure 2.9. Let the plate lie in the x - y plane, and suppose that the point charge is located at Cartesian coordinates $(0, 0, d)$. What is the scalar potential generated in the region above the plate? This is not a simple question, because the point charge induces surface charges on the plate, and we do not know beforehand how these charges are distributed.

Let us consider what do we know in this problem. We know that the conducting plate is an equipotential surface. In fact, the potential of the plate is zero, because it is grounded. We also know that the potential at infinity is zero (this is our usual boundary condition for the scalar potential). Thus, we need to solve Poisson's equation, (2.99), in the region $z > 0$, with a single point charge q located at coordinates $(0, 0, d)$, subject to the boundary conditions

$$\phi(x, y, 0) = 0, \quad (2.159)$$

and

$$\phi(x, y, z) \rightarrow 0 \quad \text{as } x^2 + y^2 + z^2 \rightarrow \infty. \quad (2.160)$$

Let us forget about the real problem, for a moment, and concentrate on a slightly different one. We shall refer to this as the *analog problem*. See Figure 2.9. In the analog problem, we have a charge q located at coordinates $(0, 0, d)$, and a charge $-q$ located at coordinates $(0, 0, -d)$, with

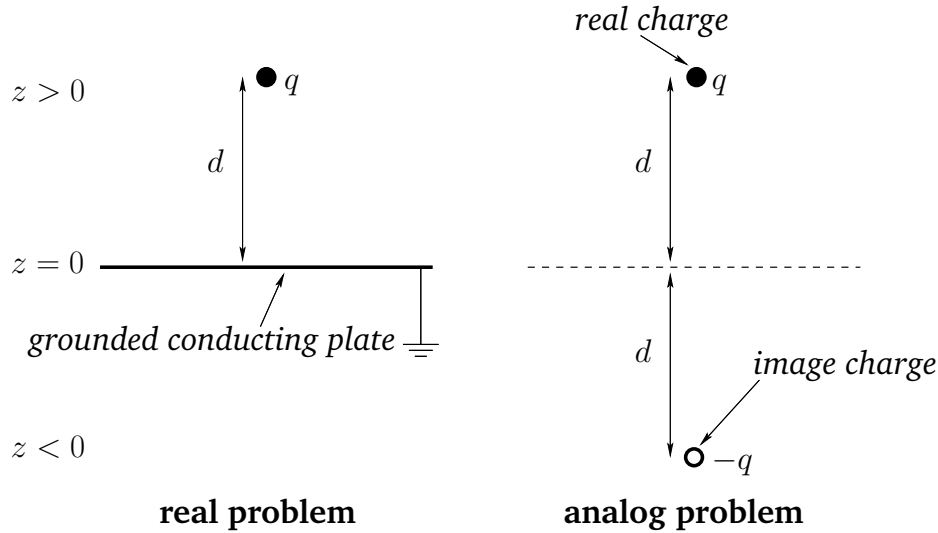


Figure 2.9: The method of images for a charge and a grounded conducting plane.

no conductors present. We can easily find the scalar potential for this problem, because we know where all the charges are located. We get

$$\phi_{\text{analog}}(x, y, z) = \frac{1}{4\pi\epsilon_0} \left[\frac{q}{\sqrt{x^2 + y^2 + (z-d)^2}} - \frac{q}{\sqrt{x^2 + y^2 + (z+d)^2}} \right]. \quad (2.161)$$

[See Equation (2.20).] Note, however, that

$$\phi_{\text{analog}}(x, y, 0) = 0, \quad (2.162)$$

and

$$\phi_{\text{analog}}(x, y, z) \rightarrow 0 \quad \text{as } x^2 + y^2 + z^2 \rightarrow \infty. \quad (2.163)$$

Moreover, in the region $z > 0$, ϕ_{analog} satisfies Poisson's equation, (2.99), for a point charge q located at coordinates $(0, 0, d)$. Thus, in this region, ϕ_{analog} is a solution to the problem posed earlier. Now, the uniqueness theorem tells us that there is only one solution to Poisson's equation that satisfies a given well-posed set of boundary conditions. (See Section 2.1.10.) So, ϕ_{analog} must be the correct potential in the region $z > 0$. Of course, ϕ_{analog} is completely wrong in the region $z < 0$. We know this because the grounded plate shields the region $z < 0$ from the point charge, so that $\phi = 0$ in this region.

Now that we have found the potential in the region $z > 0$, we can easily work out the distribution of charges induced on the conducting plate. We already know that the relation between the electric field immediately above a conducting surface and the density of charge on the surface is

$$E_{\perp} = \frac{\sigma}{\epsilon_0}. \quad (2.164)$$

[See Equation (2.129).] In this case,

$$E_{\perp}(x, y) = E_z(x, y, 0_+) = -\frac{\partial\phi(x, y, 0_+)}{\partial z} = -\frac{\partial\phi_{\text{analog}}(x, y, 0_+)}{\partial z}, \quad (2.165)$$

so

$$\sigma(x, y) = -\epsilon_0 \frac{\partial \phi_{\text{analog}}(x, y, 0_+)}{\partial z}. \quad (2.166)$$

It follows from Equation (2.161) that

$$\frac{\partial \phi_{\text{analog}}}{\partial z} = \frac{q}{4\pi \epsilon_0} \left\{ \frac{-(z-d)}{[x^2 + y^2 + (z-d)^2]^{3/2}} + \frac{(z+d)}{[x^2 + y^2 + (z+d)^2]^{3/2}} \right\}, \quad (2.167)$$

so

$$\sigma(x, y) = -\frac{q d}{2\pi (x^2 + y^2 + d^2)^{3/2}}. \quad (2.168)$$

Clearly, the charge induced on the plate has the opposite sign to the point charge. The charge density on the plate is also symmetric about the z -axis, and is largest where the plate is closest to the point charge. The total charge induced on the plate is

$$Q = \int_{x-y \text{ plane}} \sigma dS, \quad (2.169)$$

which yields

$$Q = -\frac{q d}{2\pi} \int_0^\infty \frac{2\pi r dr}{(r^2 + d^2)^{3/2}}, \quad (2.170)$$

where $r^2 = x^2 + y^2$. Thus,

$$Q = -\frac{q d}{2} \int_0^\infty \frac{dk}{(k + d^2)^{3/2}} = q d \left[\frac{1}{(k + d^2)^{1/2}} \right]_0^\infty = -q. \quad (2.171)$$

So, the total charge induced on the plate is equal and opposite to the point charge that induces it.

As we have just seen, our point electric charge induces charges of the opposite sign on the conducting plate. This, presumably, gives rise to a force of attraction between the charge and the plate. What is this force? Well, because the potentials, and, hence, the electric fields, in the vicinity of the point charge are the same in the real and analog problems, the forces acting on this charge must be the same as well. In the analog problem, there are two charges $\pm q$ a net distance $2d$ apart. The force acting on the charge at coordinates $(0, 0, d)$ (i.e., the real charge) is

$$\mathbf{f} = -\frac{q^2}{16\pi \epsilon_0 d^2} \mathbf{e}_z. \quad (2.172)$$

[See Equation (2.2).] Hence, this is also the force acting on the charge in the real problem.

What, finally, is the potential energy of the system. For the analog problem this is simply

$$W_{\text{analog}} = -\frac{q^2}{8\pi \epsilon_0 d}. \quad (2.173)$$

[See Equation (2.69).] Note that in the analog problem the fields on opposite sides of the conducting plate are mirror images of one another, as are the charges (apart from the change in sign). This is why the technique of replacing conducting surfaces by imaginary charges is called the *method*

of images. We know that the potential energy of a set of charges is equivalent to the energy stored in the surrounding electric field. Thus,

$$W = \frac{\epsilon_0}{2} \int_{\text{all space}} E^2 dV. \quad (2.174)$$

[See Equation (2.84).] Moreover, as we just mentioned, in the analog problem, the fields on either side of the x - y plane are mirror images of one another, so that $E^2(x, y, -z) = E^2(x, y, z)$. It follows that

$$W_{\text{analog}} = 2 \frac{\epsilon_0}{2} \int_{z>0} E_{\text{analog}}^2 dV. \quad (2.175)$$

Now, in the real problem,

$$\mathbf{E} = \begin{cases} \mathbf{E}_{\text{analog}} & \text{for } z > 0 \\ \mathbf{0} & \text{for } z < 0 \end{cases}. \quad (2.176)$$

So,

$$W = \frac{\epsilon_0}{2} \int_{z>0} E^2 dV = \frac{\epsilon_0}{2} \int_{z>0} E_{\text{analog}}^2 dV = \frac{1}{2} W_{\text{analog}}, \quad (2.177)$$

giving

$$W = -\frac{q^2}{16\pi \epsilon_0 d}. \quad (2.178)$$

There is another method by which we can obtain the previous result. Suppose that the point electric charge is gradually moved toward the plate along the z -axis, starting from infinity, until it reaches its final coordinates $(0, 0, d)$. How much work is required to achieve this? We know that the force of attraction acting on the charge is

$$f_z = -\frac{q^2}{16\pi \epsilon_0 z^2}. \quad (2.179)$$

[See Equation (2.172).] Thus, the work required to move this charge by dz is

$$dW = -f_z dz = \frac{q^2}{16\pi \epsilon_0 z^2} dz. \quad (2.180)$$

So, the total work needed to move the charge from $z = \infty$ to $z = d$ is

$$W = \frac{1}{4\pi \epsilon_0} \int_{\infty}^d \frac{q^2}{4z^2} dz = \frac{1}{4\pi \epsilon_0} \left[-\frac{q^2}{4z} \right]_{\infty}^d = -\frac{q^2}{16\pi \epsilon_0 d}. \quad (2.181)$$

Of course, this work is equivalent to the potential energy (2.178), and is, in turn, the same as the energy contained in the surrounding electric field.

As a second example of the method of images, consider a grounded conducting sphere of radius a centered on the origin. Suppose that a point electric charge q is placed outside the sphere at Cartesian coordinates $(b, 0, 0)$, where $b > a$. See Figure 2.10. What is the force of attraction between the sphere and the charge? In this case, we proceed by considering an analog problem in

which the sphere is replaced by an image charge $-q'$ placed somewhere on the x -axis at coordinates $(c, 0, 0)$. See Figure 2.10. The electric potential throughout space in the analog problem is simply

$$\phi(x, y, z) = \frac{q}{4\pi\epsilon_0} \frac{1}{[(x-b)^2 + y^2 + z^2]^{1/2}} - \frac{q'}{4\pi\epsilon_0} \frac{1}{[(x-c)^2 + y^2 + z^2]^{1/2}}. \quad (2.182)$$

[See Equation (2.20).] Now, the image charge must be chosen so as to make the surface $\phi = 0$ correspond to the surface of the sphere. Setting the previous expression to zero, and performing a little algebra, we find that the $\phi = 0$ surface corresponds to

$$x^2 + \frac{2(c-\lambda b)}{\lambda-1}x + y^2 + z^2 = \frac{c^2 - \lambda b^2}{\lambda-1}, \quad (2.183)$$

where $\lambda = q'^2/q^2$. Of course, the surface of the sphere satisfies

$$x^2 + y^2 + z^2 = a^2. \quad (2.184)$$

The previous two equations can be made identical by setting $\lambda = c/b$ and $a^2 = \lambda b^2$, or

$$q' = \frac{a}{b}q, \quad (2.185)$$

and

$$c = \frac{a^2}{b}. \quad (2.186)$$

According to the uniqueness theorem, the potential in the analog problem is now identical with that in the real problem in the region outside the sphere. (Of course, in the real problem, the potential inside the sphere is zero.) Hence, the force of attraction between the sphere and the original charge in the real problem is the same as the force of attraction between the image charge and the real charge in the analog problem. It follows that

$$f = \frac{qq'}{4\pi\epsilon_0(b-c)^2} = \frac{q^2}{4\pi\epsilon_0} \frac{ab}{(b^2 - a^2)^2}. \quad (2.187)$$

[See Equation (2.2).]

What is the total charge induced on the grounded conducting sphere? Well, according to Gauss's law, the flux of the electric field out of a spherical Gaussian surface lying just outside the surface of the conducting sphere is equal to the enclosed charge divided by ϵ_0 . (See Section 2.1.6.) In the real problem, the enclosed charge is the net charge induced on the surface of the sphere. In the analog problem, the enclosed charge is simply $-q'$. However, the electric fields outside the conducting sphere are identical in the real and analog problems. Hence, from Gauss's law, the charge enclosed by the Gaussian surface must also be the same in both problems. We thus conclude that the net charge induced on the surface of the conducting sphere is

$$-q' = -\frac{a}{b}q. \quad (2.188)$$

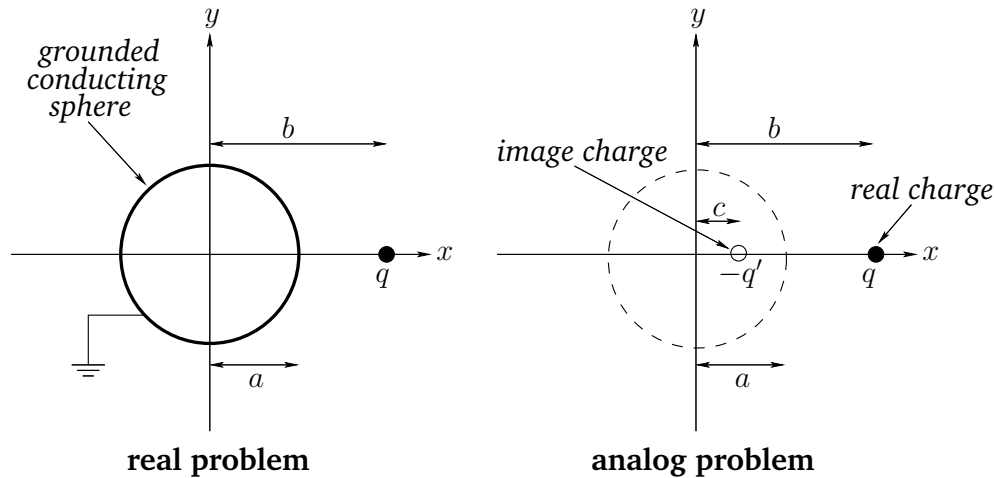


Figure 2.10: The method of images for a charge and a grounded conducting sphere.

As another example of the method of images, consider an insulated, uncharged, conducting sphere of radius a , centered on the origin, in the presence of a point electric charge q placed outside the sphere at Cartesian coordinates $(b, 0, 0)$, where $b > a$. See Figure 2.11. What is the force of attraction between the sphere and the charge? Clearly, this new problem is very similar to the one that we just discussed. The only difference is that the surface of the sphere is now at some unknown fixed potential V , and also carries zero net charge. Note that if we add a second image charge q'' , located at the origin, to the analog problem pictured in Figure 2.10 then the surface $r = a$ remains an equipotential surface. In fact, the potential of this surface becomes $V = q''/(4\pi\epsilon_0 a)$. [See Equation (2.20).] Moreover, the total charge enclosed by the surface is $-q' + q''$. This, of course, is the net charge induced on the surface of the sphere in the real problem. Hence, we can see that if $q'' = q' = (a/b)q$ then zero net charge is induced on the surface of the sphere. Thus, our modified analog problem is now a solution to the problem under discussion, in the region outside the sphere. See Figure 2.11. It follows that the surface of the sphere is at potential

$$V = \frac{q'}{4\pi\epsilon_0 a} = \frac{q}{4\pi\epsilon_0 b}. \quad (2.189)$$

Moreover, the force of attraction between the sphere and the original charge in the real problem is the same as the force of attraction between the image charges and the real charge in the analog problem. Hence, the force is given by

$$f = \frac{q q'}{4\pi\epsilon_0 (b-c)^2} - \frac{q q'}{4\pi\epsilon_0 b^2} = \frac{q^2}{4\pi\epsilon_0} \left(\frac{a}{b}\right)^3 \frac{(2b^2 - a^2)}{(b^2 - a^2)^2}. \quad (2.190)$$

[See Equation (2.2).]

As a final example of the method of images, consider two identical, infinitely long, conducting cylinders of radius a that run parallel to the z -axis, and lie a distance $2d$ apart. Suppose that one of the conductors is held at potential $+V$, while the other is held at potential $-V$. See Figure 2.12. What is the capacitance per unit length of the cylinders?

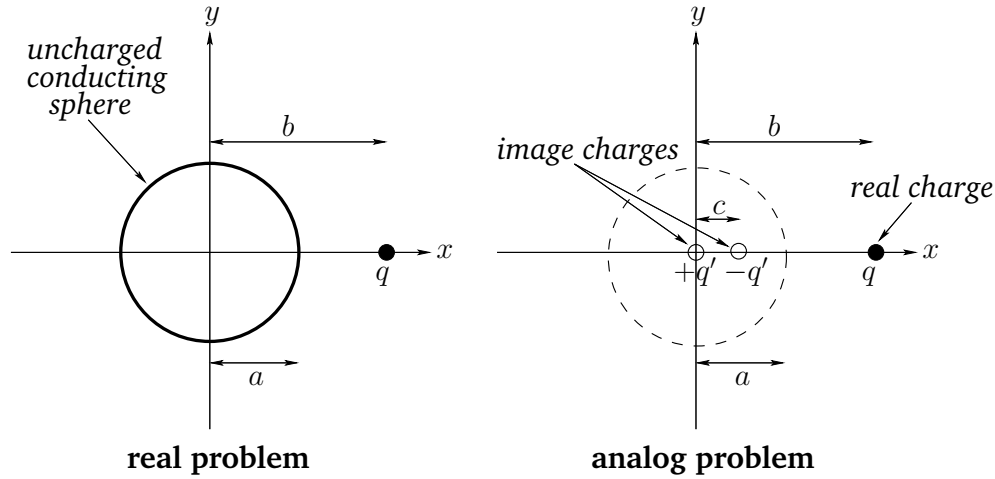


Figure 2.11: The method of images for a charge and an uncharged conducting sphere.

Consider an analog problem in which the conducting cylinders are replaced by two infinitely long charge lines, of charge per unit length $\pm\lambda$, that run parallel to the z -axis, and lie a distance $2p$ apart. Now, the potential in the x - y plane generated by a charge line λ running along the z -axis is

$$\phi(x, y) = -\frac{\lambda}{2\pi\epsilon_0} \ln r, \quad (2.191)$$

where $r = \sqrt{x^2 + y^2}$ is the radial cylindrical polar coordinate. (See Section A.23.) The corresponding electric field is radial, and satisfies

$$E_r(r) = -\frac{\partial\phi}{\partial r} = \frac{\lambda}{2\pi\epsilon_0 r}. \quad (2.192)$$

Incidentally, it is easily demonstrated from Gauss's law (see Section 2.1.6) that this is the correct electric field. Hence, the potential generated by two charge lines $\pm\lambda$ located in the x - y plane at coordinates $(\pm p, 0)$, respectively, is

$$\begin{aligned} \phi(x, y) &= \frac{\lambda}{2\pi\epsilon_0} \ln \left[\frac{1}{\sqrt{(x-p)^2 + y^2}} \right] - \frac{\lambda}{2\pi\epsilon_0} \ln \left[\frac{1}{\sqrt{(x+p)^2 + y^2}} \right] \\ &= \frac{\lambda}{4\pi\epsilon_0} \ln \left[\frac{(x+p)^2 + y^2}{(x-p)^2 + y^2} \right]. \end{aligned} \quad (2.193)$$

Suppose that

$$\frac{(x+p)^2 + y^2}{(x-p)^2 + y^2} = \alpha, \quad (2.194)$$

where α is a constant. It follows that

$$x^2 - 2p \frac{(\alpha+1)}{(\alpha-1)} x + p^2 + y^2 = 0. \quad (2.195)$$

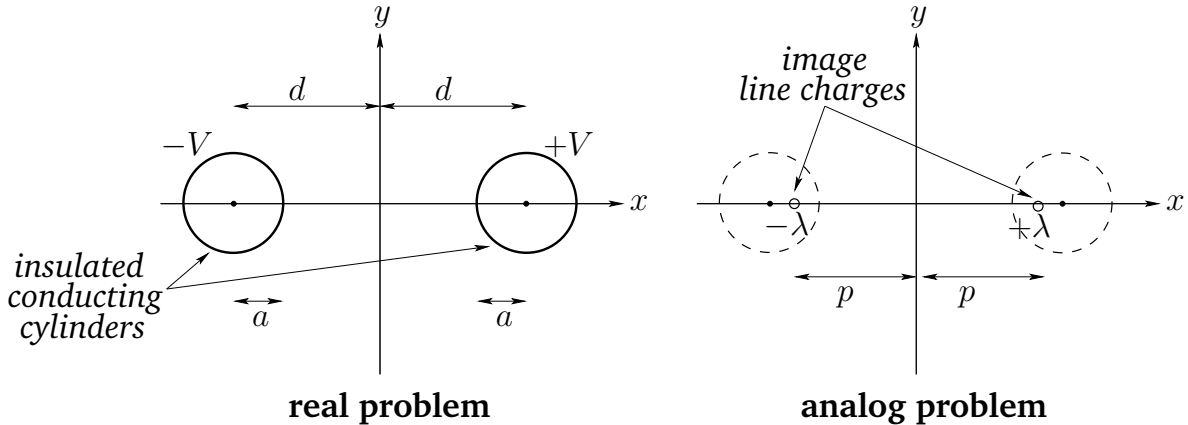


Figure 2.12: The method of images for two parallel cylindrical conductors.

Completing the square, we obtain

$$(x - d)^2 + y^2 = a^2, \quad (2.196)$$

where

$$d = \frac{(\alpha + 1)}{(\alpha - 1)} p, \quad (2.197)$$

and

$$a^2 = d^2 - p^2. \quad (2.198)$$

Of course, Equation (2.196) is the equation of a cylindrical surface of radius a centered on coordinates $(d, 0)$. Moreover, it follows from Equations (2.193) and (2.194) that this surface lies at the constant potential

$$V = \frac{\lambda}{4\pi\epsilon_0} \ln \alpha. \quad (2.199)$$

Finally, it is easily demonstrated that the equipotential $\phi = -V$ corresponds to a cylindrical surface of radius a centered on $(-d, 0)$. Hence, we can make the analog problem match the real problem in the region outside the cylinders by choosing

$$\alpha = \frac{d + p}{d - p} = \frac{d + \sqrt{d^2 - a^2}}{d - \sqrt{d^2 - a^2}}. \quad (2.200)$$

Thus, we obtain

$$V = \frac{\lambda}{4\pi\epsilon_0} \ln \left(\frac{d + \sqrt{d^2 - a^2}}{d - \sqrt{d^2 - a^2}} \right). \quad (2.201)$$

Now, it follows from Gauss's law (see Section 2.1.6), and the fact that the electric fields in the real and analog problems are identical outside the cylinders, that the charge per unit length stored on the surfaces of the two cylinders is $\pm\lambda$. Moreover, the voltage difference between the cylinders is $2V$. Hence, the capacitance per unit length of the cylinders is $C = \lambda/(2V)$, yielding

$$C = 2\pi\epsilon_0 \left/ \ln \left(\frac{d + \sqrt{d^2 - a^2}}{d - \sqrt{d^2 - a^2}} \right) \right. . \quad (2.202)$$

This expression simplifies to give

$$C = \pi \epsilon_0 \left/ \ln \left(\frac{d}{a} + \sqrt{\frac{d^2}{a^2} - 1} \right) \right., \quad (2.203)$$

which can also be written

$$C = \frac{\pi \epsilon_0}{\cosh^{-1}(d/a)}, \quad (2.204)$$

because $\cosh^{-1} x \equiv \ln(x + \sqrt{x^2 - 1})$.

2.2 Magnetostatic Fields

2.2.1 Magnetism

The phenomenon of magnetism has been known to humankind for many thousands of years. Loadstone (a magnetized form of the commonly occurring iron oxide mineral magnetite) was the first permanent magnetic material to be identified and studied. The ancient Greeks were aware of the ability of loadstone to attract small pieces of iron. The Greek word *Magnes* (*Μάγνης*), which is the root of the English word *magnet*, refers to a something (in this case, a stone) originating from Magnesia ad Sipylum, which was an ancient city in Asia Minor that was once a copious source of loadstones.

The magnetic compass was invented some time during the first ten centuries CE. Credit is variously given to the Chinese, the Arabs, and the Italians. What is certain is that, by the 12th century, magnetic compasses were in regular use by mariners to aid navigation at sea. In the 13th century, Peter Perigrinus discovered that the magnetic effect of a spherical loadstone is strongest at two oppositely directed points on the surface of the sphere, which he termed the *poles* of the magnet. He found that there are two types of poles, and that like poles repel one another, whereas unlike poles attract. In 1600, the physician William Gilbert concluded, quite correctly, that the reason that magnets preferentially align themselves in a north-south direction is that the Earth itself is a magnet. Furthermore, the Earth's magnetic poles are aligned, more or less, along its axis of rotation. This insight immediately gave rise to a fairly obvious nomenclature for the two different poles of a magnet; a magnetic *north pole* (N) has the same magnetic polarity as the geographic south pole of the Earth, and a magnetic *south pole* (S) has the same polarity as the geographic north pole of the Earth. Thus, the north pole of a magnet preferentially points northward toward the geographic north pole of the Earth (which is its magnetic south pole). In 1750, John Michell, discovered that the attractive and repulsive forces between the poles of magnets vary inversely as the square of the distance of separation. Thus, the inverse square law for forces between magnets was actually discovered prior to that for forces between electric charges.

2.2.2 Magnetic Field

In 1820, the physicist Hans Christian Ørsted was giving a lecture demonstration of various electrical and magnetic effects. Suddenly, much to his amazement, he noticed that the needle of a

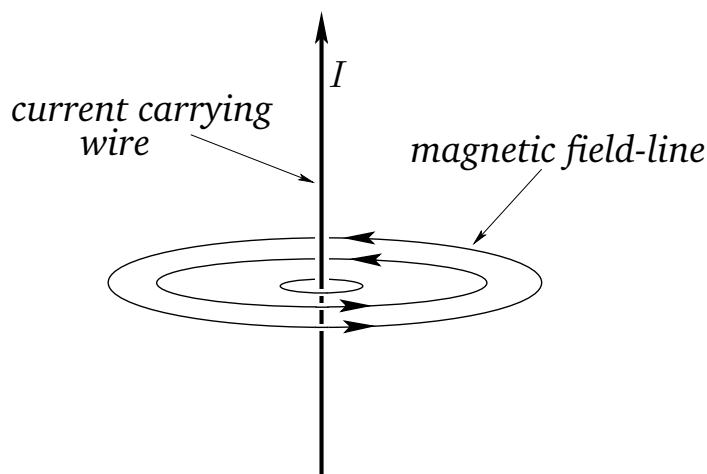


Figure 2.13: Magnetic loops around a current-carrying wire.

compass that he was holding was deflected when he moved it close to a current-carrying wire. This was a very surprising observation, because, until that moment, electricity and magnetism had been thought of as two quite unrelated phenomena. Word of this discovery spread quickly, and scientists such as Andre Marie Ampère, François Arago, Jean-Baptiste Biot, Félix Savart, and Michael Faraday immediately decided to investigate further. Their discoveries can be encapsulated by describing a series of simple, and easily reproducible, experiments.

Consider an experiment in which a long straight wire carries an electrical current I . As is easily demonstrated, the needle of a small compass maps out a series of concentric circular loops in the plane perpendicular to such a wire. See Figure 2.13. The direction of circulation around such magnetic loops is conventionally taken to be the direction in which the north pole of a compass needle points. Using this convention, the circulation of the loops is given by a *right-hand rule*. If the thumb of the right-hand points along the direction of the current then the fingers of the right-hand circulate in the same sense as the magnetic loops.

Our next experiment involves bringing a short test wire, carrying a current I' , close to the original long straight wire, and investigating the force exerted on the test wire. This experiment is not quite as clear cut as Coulomb's experiment regarding the force exerted between electric charges, because, unlike electric charges, electric currents cannot exist as point entities; they have to flow in complete circuits. We must imagine that the circuit that connects with the central wire is sufficiently far away that it has no appreciable influence on the outcome of the experiment. The circuit that connects with the test wire is more problematic. Fortunately, if the feed wires are twisted around each other, as indicated in Figure 2.14, then they effectively cancel one another out, and also do not influence the outcome of the experiment.

It can easily be demonstrated that the force exerted on the test wire is directly proportional to its length. Furthermore, if the current in the test wire (i.e., the test current) flows parallel to the current in the central wire then the two wires attract one another. If the current in the test wire is reversed then the two wires repel one another. If the test current is directed radially toward the central wire (and the current in the central wire flows upward) then the test wire is subject to a downward

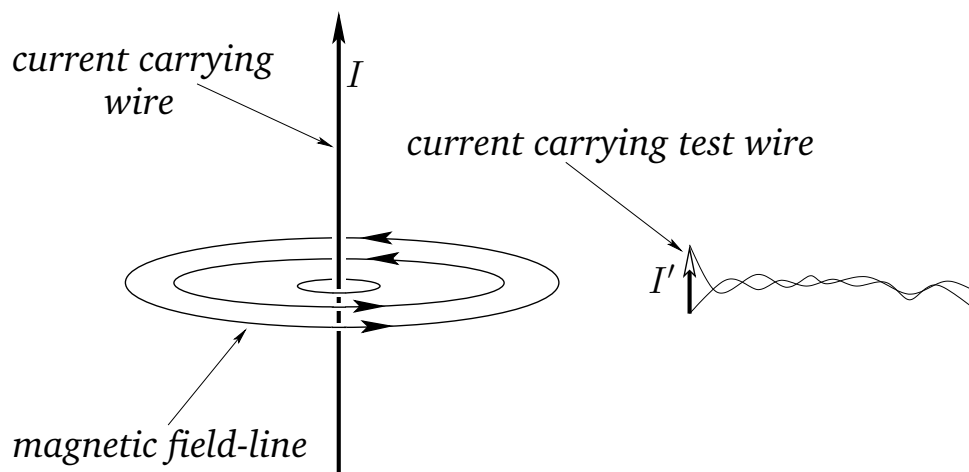


Figure 2.14: Force on current-carrying wire.

force. If the test current is reversed then the force is upward. If the test current is rotated in a single plane, such that it starts parallel to the central current and ends up pointing radially toward it, then the force on the test wire is of constant magnitude, and is always perpendicular to the test current. If the test current is parallel to a magnetic loop then there is no force exerted on the test wire. If the test current is rotated in a single plane, such that it starts parallel to the central current, and ends up pointing along a magnetic loop, then the magnitude of the force on the test wire attenuates like $\cos \theta$ (where θ is the angle through which the current is turned, and $\theta = 0$ corresponds to the case where the test current is parallel to the central current), and its direction is again always perpendicular to the test current. Finally, the attractive force between two parallel current-carrying wires is proportional to the product of the two currents, and inversely proportional to the perpendicular distance between the wires.

The rather complicated force law established by the previously described experiments can be summed up succinctly provided that we define a vector field $\mathbf{B}(\mathbf{r})$, called the *magnetic field*, that fills space, and whose direction is everywhere tangential to the magnetic loops mapped out by the north pole of a small compass. The dependence of the force per unit length, \mathbf{F} , acting on a test wire, located at displacement \mathbf{r} , with the different possible orientations of the test current is described by

$$\mathbf{F}(\mathbf{r}) = \mathbf{I}' \times \mathbf{B}(\mathbf{r}), \quad (2.205)$$

where \mathbf{I}' is a vector whose direction and magnitude are the same as those of the test current.

The variation of the force per unit length acting on a test wire with the strength of the central current, I , and the perpendicular distance, r , to the central wire, is accounted for by saying that the strength of the magnetic generated around the central wire is directly proportional to I , and inversely proportional to r . Thus, we can write

$$B = \frac{\mu_0 I}{2\pi r}. \quad (2.206)$$

The constant of proportionality μ_0 is called the *magnetic permeability of free space*, and takes the

value

$$\mu_0 = 4\pi \times 10^{-7} \text{ N A}^{-2}. \quad (2.207)$$

Incidentally, the SI unit of magnetic field strength is the *tesla* (T), which is equivalent to a newton per ampere per meter.

The concept of a magnetic field that fills the space around a current-carrying wire allows the calculation of the force on a test wire to be conveniently split into two parts. In the first part, we calculate the magnetic field generated by the current flowing in the central wire. This field circulates in the plane normal to the wire. Its magnitude is proportional to the central current, and inversely proportional to the perpendicular distance from the wire. In the second part, we employ Equation (2.205) to calculate the force per unit length acting on a short current-carrying wire placed in the magnetic field generated by the central current. This force is perpendicular to both the direction of the magnetic field and the direction of the test current. Note that, at this stage, we have no reason to suppose that the magnetic field has any real existence; it is introduced merely to facilitate the calculation of the force exerted on the test wire by the central wire.

2.2.3 Ampère's Law

It is an experimentally demonstrable fact that magnetic fields, like electric fields, are completely superposable. So, if a magnetic field $\mathbf{B}_1(\mathbf{r})$ is generated by an electric current I_1 flowing through some circuit, and a field $\mathbf{B}_2(\mathbf{r})$ is generated by a current I_2 flowing through another circuit, then when the currents I_1 and I_2 flow through both circuits simultaneously the generated magnetic field is $\mathbf{B}_1(\mathbf{r}) + \mathbf{B}_2(\mathbf{r})$.

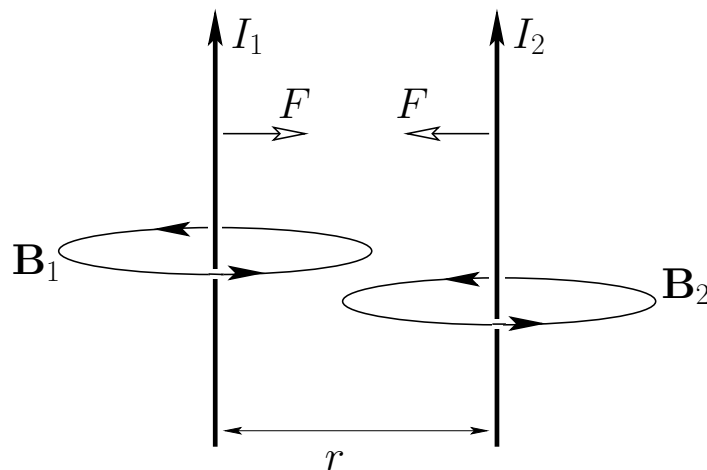


Figure 2.15: Two parallel current-carrying wires.

Consider two parallel wires separated by a perpendicular distance r , and carrying electric currents I_1 and I_2 , respectively. The magnetic field-strength at the second wire due to the current flowing in the first wire is $B = \mu_0 I_1 / 2\pi r$. [See Equation (2.206).] This field is orientated perpen-

dicular to the second wire, so the force per unit length exerted on the second wire is

$$F = \frac{\mu_0 I_1 I_2}{2\pi r}. \quad (2.208)$$

The previous expression follows from Equation (2.205), which is valid for continuous wires as well as short test wires. The force acting on the second wire is directed radially inward toward the first wire. The magnetic field-strength at the first wire due to the current flowing in the second wire is $B = \mu_0 I_2/2\pi r$. This field is orientated perpendicular to the first wire, so the force per unit length acting on the first wire is equal and opposite to that acting on the second wire, according to Equation (2.205). Equation (2.208) is known as *Ampère's law*.

Equation (2.208) is the basis of the official (prior to 2019) SI definition of the ampere, which is:

One ampere is the magnitude of the current which, when flowing in each of two long parallel wires one meter apart, results in a force between the wires of 2×10^{-7} N per meter of length.

We can see that it is no accident that the constant μ_0 has the numerical value of exactly $4\pi \times 10^{-7}$. (Incidentally, this rather strange definition arose because electromagnetism was originally formulated in the cgs system of units. In the cgs system, the force per unit length exerted by two parallel wires, one centimeter apart, both carrying a current of 1 abampere (i.e., 10 amperes), is 2 dynes per centimeter.)

2.2.4 Lorentz Force Law

The flow of an electric current down a conducting wire is ultimately due to the movement of electrically charged particles (in most cases, electrons) along the wire. It seems reasonable, therefore, that the force exerted on the wire when it is placed in a magnetic field is simply the resultant of the forces exerted on these moving charges. Let us suppose that this is the case.

Let A be the (uniform) cross-sectional area of the wire, and let n be the number density of mobile charges in the wire. Suppose that the mobile charges each have charge q and drift velocity \mathbf{v} . We must assume that the wire also contains stationary charges, of charge $-q$ and number density n , say, so that the net charge density in the wire is zero. In most conductors, the mobile charges are electrons, and the stationary charges are ions. The magnitude of the electric current flowing through the wire is simply the number of coulombs per second that flow past a given point. In one second, a mobile charge moves a distance v , so all of the charges contained in a cylinder of cross-sectional area A and length v flow past a given point. Thus, the magnitude of the current is $qnAv$. The direction of the current is the same as the direction of motion of the charges (i.e., $\mathbf{I} \propto \mathbf{v}$), so the vector current is

$$\mathbf{I}' = qnA\mathbf{v}. \quad (2.209)$$

According to Equation (2.205), the force per unit length acting on the wire is

$$\mathbf{F} = \mathbf{I}' \times \mathbf{B} = qnA\mathbf{v} \times \mathbf{B}. \quad (2.210)$$

However, a unit length of the wire contains nA moving charges. So, assuming that each charge is subject to an equal force from the magnetic field (and we have no reason to suppose otherwise), the magnetic force acting on an individual charge is

$$\mathbf{f} = q \mathbf{v} \times \mathbf{B}. \quad (2.211)$$

This formula implies that the magnitude of the magnetic force exerted on a moving charged particle is the product of the particle's electric charge, its velocity, the magnetic field-strength, and the sine of the angle subtended between the particle's direction of motion and the direction of the magnetic field. (See Section A.8.) The force is directed perpendicular to both the magnetic field and the particle's instantaneous direction of motion.

We can combine the previous equation with Equation (2.10) to give the force acting on an electric charge q moving with velocity \mathbf{v} in an electric field \mathbf{E} and a magnetic field \mathbf{B} :

$$\mathbf{f} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}). \quad (2.212)$$

This result is called the *Lorentz force law*, after Hendrick Antoon Lorentz, who first formulated it. The electric force on a charged particle is parallel to the local electric field. The magnetic force, however, is perpendicular to both the local magnetic field and the particle's direction of motion. No magnetic force is exerted on a stationary charged particle.

The equation of motion of a free particle of charge q and mass m moving in electric and magnetic fields is

$$m \mathbf{a} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}), \quad (2.213)$$

according to the Lorentz force law. (See Section 1.2.3.) Here, \mathbf{a} is the particle's acceleration. This equation of motion was verified in a famous experiment carried out by the Cambridge physicist J.J. Thompson in 1897. Thompson was investigating *cathode rays*, a then mysterious form of radiation emitted by a heated metal element held at a large negative voltage (i.e., a cathode) with respect to another metal element (i.e., an anode) in an evacuated tube. German physicists maintained that cathode rays were a form of electromagnetic radiation, whereas British and French physicists suspected that they were, in reality, a stream of charged particles. Thompson was able to demonstrate that the latter view was correct. In Thompson's experiment, the cathode rays pass through a region of crossed electric and magnetic fields (still in vacuum). The fields are perpendicular to the original trajectory of the rays, and are also mutually perpendicular.

Let us analyze Thompson's experiment. Suppose that the rays are originally traveling in the x -direction, and are subject to a uniform electric field E in the z -direction, and a uniform magnetic field B in the $-y$ -direction. See Figure 2.16. Let us assume, as Thompson did, that cathode rays are a stream of particles of mass m and charge q . The z -component of the equation of motion of an individual particle is

$$m a_z = q(E - vB), \quad (2.214)$$

where v is the x -component of its velocity, and a_z the z -component of its acceleration. Thompson started off his experiment by only turning on the electric field in his apparatus, and measuring the deflection d of the rays in the z -direction after they had traveled a distance l through the field.

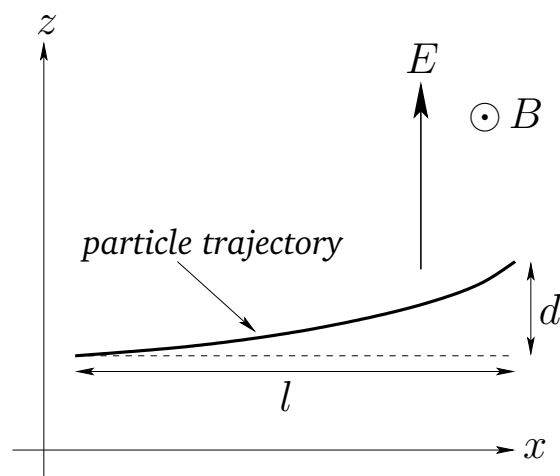


Figure 2.16: Thomson's experiment.

Now, a particle subject to a constant acceleration a_z in the z -direction is deflected a distance $d = (1/2) a_z t^2$ in a time t . Thus,

$$d = \frac{1}{2} \frac{qE}{m} t^2 = \frac{qE l^2}{m 2v^2}, \quad (2.215)$$

where the time of flight t is replaced by l/v . This replacement is only valid if $d \ll l$ (i.e., if the deflection of the rays is small compared to the distance that they travel through the electric field), which is assumed to be the case. Next, Thomson turned on the magnetic field in his apparatus, and adjusted it so that the cathode rays were no longer deflected. The lack of deflection implies that the net force on the particles in the z -direction is zero. In other words, the electric and magnetic forces balance exactly. It follows from Equation (2.214) that, with a properly adjusted magnetic field-strength,

$$v = \frac{E}{B}. \quad (2.216)$$

Thus, Equations (2.215) and (2.216) can be combined and rearranged to give the charge to mass ratio of the particles in terms of measured quantities:

$$\frac{q}{m} = \frac{2dE}{l^2 B^2}. \quad (2.217)$$

Using this method, Thomson inferred that cathode rays are made up of negatively charged particles (the sign of the charge is obvious from the direction of the deflection in the electric field) with a charge to mass ratio of $-1.7 \times 10^{11} \text{ C kg}^{-1}$.

A decade later, in 1908, Robert Millikan performed his famous oil drop experiment in which he discovered that mobile electric charges are quantized in units of $-1.6 \times 10^{-19} \text{ C}$. Assuming that mobile electric charges and the particles that make up cathode rays are one and the same thing, Thomson's and Millikan's experiments imply that the mass of these particles is $9.4 \times 10^{-31} \text{ kg}$. Of course, this is the mass of an electron (the modern value is $9.1 \times 10^{-31} \text{ kg}$), and $-1.6 \times 10^{-19} \text{ C}$ is the charge of an electron. Thus, cathode rays are, in fact, streams of electrons that are emitted

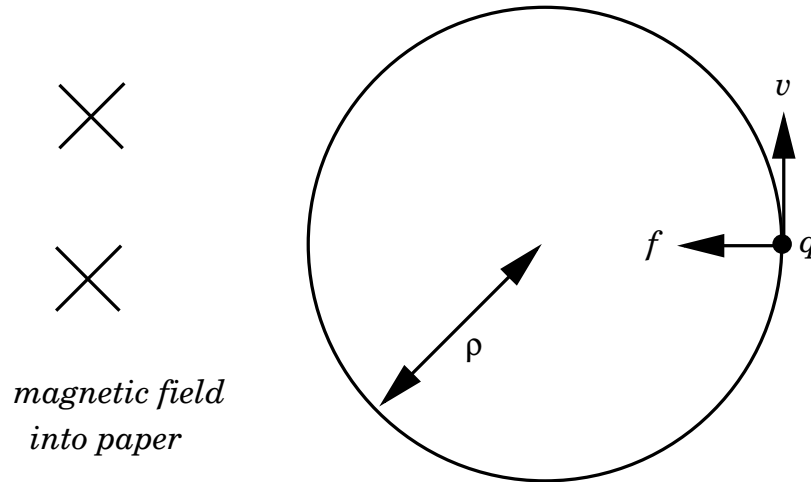


Figure 2.17: Circular motion of a charged particle in a magnetic field.

from a heated cathode, and then accelerated because of the large voltage difference between the cathode and anode.

If a particle is subject to a force \mathbf{f} that causes it to displace by $d\mathbf{r}$ then the work done on the particle by the force is

$$W = \mathbf{f} \cdot d\mathbf{r} = f dr \cos \theta, \quad (2.218)$$

where θ is the angle subtended between the force and the displacement. (See Section 1.3.2.) However, this angle is always 90° for the force exerted by a magnetic field on a charged particle, because the magnetic force is always perpendicular to the particle's instantaneous direction of motion. It follows that a magnetic field is unable to do work on a charged particle. In other words, a charged particle can never gain or lose energy due to interaction with a magnetic field. On the other hand, a charged particle can certainly gain or lose energy due to interaction with an electric field. Thus, magnetic fields are often used in particle accelerators to guide charged particle motion (e.g., in a circle), but the actual acceleration is always performed by electric fields.

2.2.5 Charged Particle Motion in a Magnetic Field

Suppose that a particle of mass m moves in a circular orbit of radius ρ with a constant speed v . As is well known, the acceleration of the particle is of magnitude v^2/ρ , and is always directed toward the center of the orbit. It follows that the acceleration is always perpendicular to the particle's instantaneous direction of motion.

We have seen that the force exerted on an electrically charged particle by a magnetic field is always perpendicular to its instantaneous direction of motion. Does this imply that the field causes the particle to execute a circular orbit? Consider the case shown in Figure 2.17. Suppose that a particle of positive charge q and mass m moves in a plane perpendicular to a uniform magnetic field B . In the figure, the field is directed into the plane of the paper. Suppose that the particle moves, in a counter-clockwise manner, with constant speed v (recall that the magnetic field cannot

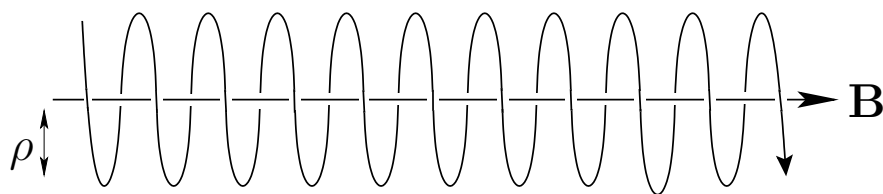


Figure 2.18: Spiral trajectory of a charged particle in a uniform magnetic field.

do work on the particle, so it cannot affect its speed), in a circular orbit of radius ρ . The magnetic force acting on the particle is of magnitude $f = qvB$ and, according to Equation (2.211), this force is always directed toward the center of the orbit. Thus, if

$$f = qvB = \frac{mv^2}{\rho}, \quad (2.219)$$

then we have a self-consistent picture. It follows that

$$\rho = \frac{mv}{qB}. \quad (2.220)$$

The angular frequency of rotation of the particle (i.e., the number of radians the particle rotates through in one second) is

$$\omega = \frac{v}{\rho} = \frac{qB}{m}. \quad (2.221)$$

Note that this frequency, which is known as the *Larmor frequency*, does not depend on the velocity of the particle. For a negatively charged particle, the picture is exactly the same as described previously, except that the particle moves in a clockwise orbit.

It is clear, from Equation (2.221), that the angular frequency of gyration of a charged particle in a known magnetic field can be used to determine its charge to mass ratio, q/m . Furthermore, if the speed of the particle is known then the radius of the orbit can also be used to determine q/m , via Equation (2.220). In the past, this method was used extensively in high energy physics experiments to identify particles from photographs of the tracks that they left in magnetized cloud chambers or bubble chambers. It is, of course, easy to differentiate positively charged particles from negatively charged ones using the direction of deflection of the particles in the magnetic field.

We have seen that a charged particle placed in a magnetic field executes a circular orbit in the plane perpendicular to the direction of the field. However, we can also add an arbitrary drift along the direction of the magnetic field. This follows because the force $q\mathbf{v} \times \mathbf{B}$ acting on the particle only depends on the component of the particle's velocity that is perpendicular to the direction of magnetic field (the vector product of two parallel vectors is always zero because the angle θ they subtend is zero). (See Section A.8.) The combination of circular motion in the plane perpendicular to the magnetic field, and uniform motion along the direction of the field, gives rise to a spiral trajectory of a charged particle in a magnetic field, where the field forms the axis of the spiral. See Figure 2.18.

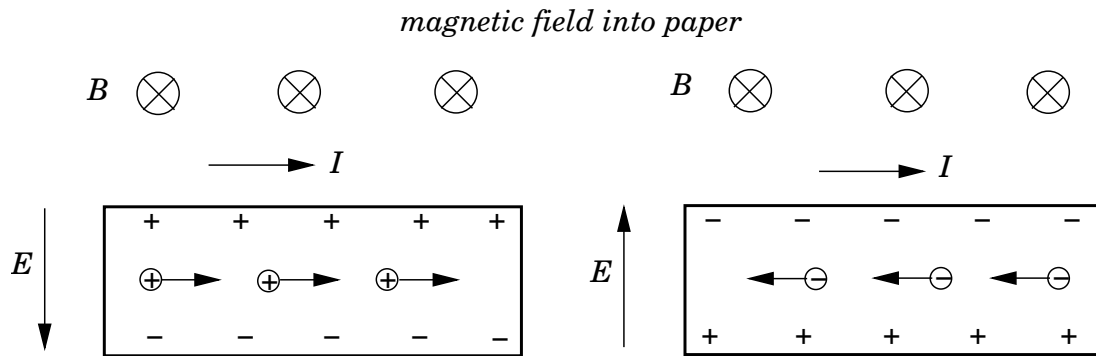


Figure 2.19: Hall effect for positive charge carriers (left) and negative charge carriers (right).

2.2.6 Hall Effect

We have repeatedly stated that the mobile charges in conventional conducting materials are negatively charged. (They are, in fact, electrons.) Is there any direct experimental evidence that this is true? Actually, there is. We can use a phenomenon called the *Hall effect* to determine whether the mobile charges in a given conductor are positively or negatively charged. Let us investigate this effect.

Consider a thin, flat, uniform, ribbon of some conducting material that is orientated such that its flat side is perpendicular to a uniform magnetic field B . See Figure 2.19. Suppose that we pass a current I along the length of the ribbon. There are two alternatives. Either the current is carried by positive charges moving from left to right (in the figure), or it is carried by negative charges moving in the opposite direction.

Suppose that the current is carried by positive charges moving from left to right. These charges are deflected upward (in the figure) by the magnetic field. Thus, the upper edge of the ribbon becomes positively charged, while the lower edge becomes negatively charged. Consequently, there is a positive potential difference V_H between the upper and lower edges of the ribbon. This potential difference is called the *Hall voltage*.

Suppose, now, that the current is carried by negative charges moving from right to left. These charges are also deflected upward by the magnetic field. Thus, the upper edge of the ribbon becomes negatively charged, while the lower edge becomes positively charged. It follows that the Hall voltage (i.e., the potential difference between the upper and lower edges of the ribbon) is negative in this case.

Clearly, it is possible to determine the sign of the mobile charges in a current-carrying conductor by measuring the Hall voltage. If the voltage is positive then the mobile charges are positive (assuming that the magnetic field and the current are orientated as shown in the figure), whereas if the voltage is negative then the mobile charges are negative. If we were to perform this experiment then we would discover that the mobile charges in metals are always negative (because they are electrons). However, in some types of semiconductor the mobile charges turn out to be positive. These positive charge carriers are called *holes*. Holes are actually missing electrons in the atomic lattice of the semiconductor, but they act essentially like positive charges.

Let us investigate the magnitude of the Hall voltage. Suppose that the mobile charges each possess a charge q and move along the ribbon with the drift velocity v_d . The magnetic force on a given mobile charge is of magnitude $q v_d B$, because the charge moves essentially perpendicular to the magnetic field. [See Equation (2.211).] In a steady state, this force is balanced by the electric force due to the build up of charges on the upper and lower edges of the ribbon. If the Hall voltage is V_H , and the width of the ribbon is w , then the electric field directed from the upper to the lower edge of the ribbon is of magnitude $E = V_H/w$. [See Equation (2.17).] Now, the electric force on a mobile charge is $q E$. [See Equation (2.10).] This force acts in opposition to the magnetic force. In a steady state,

$$q E = \frac{q V_H}{w} = q v_d B, \quad (2.222)$$

giving

$$V_H = v_d w B. \quad (2.223)$$

Note that the Hall voltage is directly proportional to the magnitude of the magnetic field. In fact, this property of the Hall voltage is exploited in instruments, called *Hall probes*, that are used to measure magnetic field-strengths.

Suppose that the thickness of the conducting ribbon is d , and that it contains n mobile charge carriers per unit volume. It follows that the total current flowing through the ribbon can be written

$$I = q n w d v_d, \quad (2.224)$$

because all mobile charges contained in a rectangular volume of length v_d , width w , and thickness d , flow past a given point on the ribbon in one second. Combining Equations (2.223) and (2.224), we obtain

$$V_H = \frac{I B}{q n d}. \quad (2.225)$$

It is clear that the Hall voltage is proportional to the current flowing through the ribbon and the magnetic field-strength, and is inversely proportional to the number density of mobile charges in the ribbon and the thickness of the ribbon. Thus, in order to construct a sensitive Hall probe (i.e., one that produces a large Hall voltage in the presence of a small magnetic field), we need to take a thin ribbon of some material that possesses relatively few mobile charges per unit volume (e.g., a semiconductor), and then run a large current through it.

2.2.7 Biot-Savart Law

Consider a closed electric circuit of general shape, fabricated from an idealized zero thickness wire, around which a current I flows. According to *Biot-Savart law*, which is named after Jean-Baptiste Biot and Félix Savart, and which can be experimentally verified, the magnetic field generated by such a circuit is

$$\mathbf{B}(\mathbf{r}) = \frac{\mu_0 I}{4\pi} \oint \frac{d\mathbf{r}' \times (\mathbf{r} - \mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|^3}, \quad (2.226)$$

where $d\mathbf{r}'$ is an element of the wire, whose displacement is \mathbf{r}' , and the integral is taken around the whole circuit.

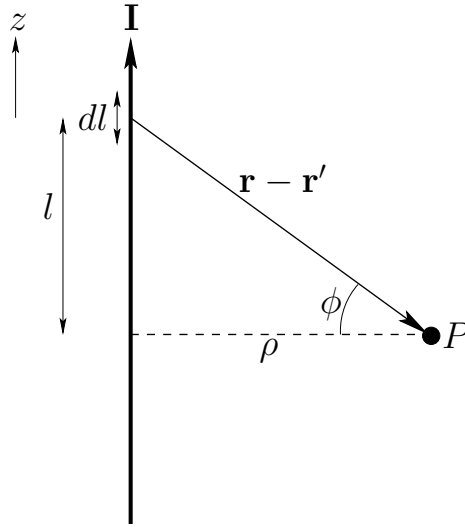


Figure 2.20: A Biot-Savart law calculation.

Consider an infinite straight wire, running along the z -axis, that carries a current I . See Figure 2.20. Let us reconstruct the magnetic field generated by the wire at point P using the Biot-Savart law. Suppose that the perpendicular distance to the wire is ρ . It is easily seen that

$$\mathbf{e}_z \times (\mathbf{r} - \mathbf{r}') = \rho \mathbf{e}_\theta, \quad (2.227)$$

$$l = \rho \tan \phi, \quad (2.228)$$

$$dl = \frac{\rho}{\cos^2 \phi} d\phi, \quad (2.229)$$

$$|\mathbf{r} - \mathbf{r}'| = \frac{\rho}{\cos \phi}, \quad (2.230)$$

where θ is a cylindrical polar coordinate. (See Section A.23.) Hence,

$$d\mathbf{r}' \times (\mathbf{r} - \mathbf{r}') = \frac{\rho^2 \mathbf{e}_\theta}{\cos^2 \phi} d\phi. \quad (2.231)$$

Thus, according to Equation (2.226), we have

$$\begin{aligned} \mathbf{B} &= \frac{\mu_0 I}{4\pi} \int_{-\pi/2}^{\pi/2} \frac{\rho^2}{\cos^2 \phi} \frac{1}{\rho^3 (\cos \phi)^{-3}} d\phi \mathbf{e}_\theta \\ &= \frac{\mu_0 I}{4\pi \rho} \int_{-\pi/2}^{\pi/2} \cos \phi d\phi \mathbf{e}_\theta = \frac{\mu_0 I}{4\pi \rho} [\sin \phi]_{-\pi/2}^{\pi/2} \mathbf{e}_\theta, \end{aligned} \quad (2.232)$$

which gives

$$\mathbf{B} = \frac{\mu_0 I}{2\pi \rho} \mathbf{e}_\theta. \quad (2.233)$$

Thus, we conclude that the Biot-Savart law is a more general form of the familiar result (2.206) that is not restricted to long straight wires.

Consider a circular wire loop of radius a that carries a current I . Suppose that the loop lies in the x - y plane, and is centered on the origin. Let us use the Biot-Savart law to calculate the magnetic field generated by the coil along a perpendicular axis that passes through its center (i.e., along the z -axis). Let z be the distance of the point of observation from the center of the loop, and let the angle θ parameterize position on the loop. Thus, we have

$$\mathbf{r} = (0, 0, z), \quad (2.234)$$

$$\mathbf{r}' = (a \cos \theta, a \sin \theta, 0), \quad (2.235)$$

where the right-hand sides of the previous two equations are Cartesian components. It follows that

$$\mathbf{r} - \mathbf{r}' = (-a \cos \theta, -a \sin \theta, z), \quad (2.236)$$

$$|\mathbf{r} - \mathbf{r}'| = (a^2 + z^2)^{1/2}, \quad (2.237)$$

$$d\mathbf{r}' = (-a \sin \theta d\theta, a \cos \theta d\theta, 0), \quad (2.238)$$

$$d\mathbf{r}' \times (\mathbf{r} - \mathbf{r}') = (a z \cos \theta d\theta, a z \sin \theta d\theta, a^2 d\theta). \quad (2.239)$$

Thus, the Biot-Savart law, (2.226), yields

$$B_x = \frac{\mu_0 I}{4\pi} \oint \frac{a z \cos \theta d\theta}{(a^2 + z^2)^{3/2}} = 0, \quad (2.240)$$

$$B_y = \frac{\mu_0 I}{4\pi} \oint \frac{a z \sin \theta d\theta}{(a^2 + z^2)^{3/2}} = 0, \quad (2.241)$$

$$B_z = \frac{\mu_0 I}{4\pi} \oint \frac{a^2 d\theta}{(a^2 + z^2)^{3/2}} = \frac{\mu_0 I}{2} \frac{a^2}{(a^2 + z^2)^{3/2}}. \quad (2.242)$$

Thus, the magnetic field generated on the z -axis is

$$\mathbf{B} = \frac{\mu_0 I}{2} \frac{a^2}{(a^2 + z^2)^{3/2}} \mathbf{e}_z. \quad (2.243)$$

Suppose that we have two identical current loops of radius a . Let both loops be centered on the z -axis, and let the first lie in the plane $z = d$, and the second in the plane $z = -d$. Furthermore, suppose that a current I flows around each loop in the same direction. By the principle of superposition, making use of the previous equation, the magnetic field generated on the z -axis by the two loops is

$$B_z = \frac{\mu_0 I}{2} \left(\frac{a^2}{[a^2 + (z - d)^2]^{3/2}} + \frac{a^2}{[a^2 + (z + d)^2]^{3/2}} \right). \quad (2.244)$$

If we Taylor expand the previous expression about $z = 0$ then we obtain

$$B_z = \frac{\mu_0 I}{2} \frac{a^2}{(a^2 + d^2)^{3/2}} \left\{ 2 + 3 \left[\frac{(2d)^2 - a^2}{(a^2 + d^2)^2} \right] z^2 + \mathcal{O}(z^4) \right\}. \quad (2.245)$$

Suppose that we wish to make the magnetic field in the region between the loops as uniform as possible. We can clearly achieve this goal if we adjust the spacing $2d$ between the loops in such a manner that the coefficient of z^2 in the previous expression is set to zero. In this case, the leading order non-constant term in the expansion is $O(z^4)$. It can be seen that we need $2d = a$. In other words, the spacing between the loops must equal the radius of the loops. The approximately uniform magnetic field between the loops becomes

$$B_z = \left(\frac{4}{5}\right)^{3/2} \frac{\mu_0 I}{a}. \quad (2.246)$$

A pair of current loops set up in this manner are known as *Helmholtz coils*.

Finally, we can generalize the Biot-Savart law, (2.226), to determine the magnetic field generated by a distributed current of density $\mathbf{j}(\mathbf{r})$ by making the identification

$$I d\mathbf{r} = \mathbf{j}(\mathbf{r}) dV. \quad (2.247)$$

Thus, we obtain

$$\mathbf{B}(\mathbf{r}) = \frac{\mu_0}{4\pi} \int \frac{\mathbf{j}(\mathbf{r}') \times (\mathbf{r} - \mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|^3} dV', \quad (2.248)$$

where the volume integral is taken over all space.

2.2.8 Magnetic Vector Potential

We saw in Equation (2.16) that

$$\frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3} = -\nabla \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right). \quad (2.249)$$

This equation can be combined with the generalized Biot-Savart law, (2.248), to give

$$\mathbf{B}(\mathbf{r}) = \frac{\mu_0}{4\pi} \int \nabla \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right) \times \mathbf{j}(\mathbf{r}') dV'. \quad (2.250)$$

It follows that

$$\mathbf{B} = \nabla \times \mathbf{A}, \quad (2.251)$$

where

$$\mathbf{A}(\mathbf{r}) = \frac{\mu_0}{4\pi} \int \frac{\mathbf{j}(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} dV'. \quad (2.252)$$

(See Section A.24.) Here, the vector field $\mathbf{A}(\mathbf{r})$ is known as the *magnetic vector potential*.

It is possible to prove that the magnetic vector potential defined in the previous equation is a divergence-free field. Note that

$$\frac{\partial}{\partial x} \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right) = -\frac{x - x'}{|\mathbf{r} - \mathbf{r}'|^3} = \frac{x' - x}{|\mathbf{r} - \mathbf{r}'|^3} = -\frac{\partial}{\partial x'} \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right), \quad (2.253)$$

which implies that

$$\nabla \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right) = -\nabla' \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right), \quad (2.254)$$

where ∇' is the operator $(\partial/\partial x', \partial/\partial y', \partial/\partial z')$. (See Section A.19.) Taking the divergence of Equation (2.252), and making use of the previous relation, we obtain

$$\nabla \cdot \mathbf{A} = \frac{\mu_0}{4\pi} \int \mathbf{j}(\mathbf{r}') \cdot \nabla \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right) dV' = -\frac{\mu_0}{4\pi} \int \mathbf{j}(\mathbf{r}') \cdot \nabla' \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right) dV'. \quad (2.255)$$

Now,

$$\int_{-\infty}^{\infty} g \frac{\partial f}{\partial x} dx = [gf]_{-\infty}^{\infty} - \int_{-\infty}^{\infty} f \frac{\partial g}{\partial x} dx. \quad (2.256)$$

However, if $gf \rightarrow 0$ as $x \rightarrow \pm\infty$ then we can neglect the first term on the right-hand side of the previous equation, and write

$$\int_{-\infty}^{\infty} g \frac{\partial f}{\partial x} dx = - \int_{-\infty}^{\infty} f \frac{\partial g}{\partial x} dx. \quad (2.257)$$

A simple generalization of this result yields

$$\int \mathbf{g} \cdot \nabla f dV = - \int f \nabla \cdot \mathbf{g} dV, \quad (2.258)$$

provided that $g_x f \rightarrow 0$ as $|\mathbf{r}| \rightarrow \infty$, etc cetera. Thus, Equation (2.255) yields

$$\nabla \cdot \mathbf{A} = \frac{\mu_0}{4\pi} \int \frac{\nabla' \cdot \mathbf{j}(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} dV', \quad (2.259)$$

provided that $|\mathbf{j}(\mathbf{r})|$ is bounded as $|\mathbf{r}| \rightarrow \infty$. Now, the flux of electric charge out of a surface S , enclosing a volume V , is

$$\oint_S \mathbf{j} \cdot d\mathbf{S} = \int_V \nabla \cdot \mathbf{j} dV, \quad (2.260)$$

where use has been made of the divergence theorem. (See Section A.20.) However, for a steady current distribution, this flux must be zero, otherwise positive or negative electric charge would build up inside V . Moreover, the flux must be zero for all possible volumes, V , which implies that

$$\nabla \cdot \mathbf{j} = 0 \quad (2.261)$$

for a steady current distribution. Hence, we deduce from Equation (2.259) that

$$\nabla \cdot \mathbf{A} = 0 \quad (2.262)$$

for a steady current distribution.

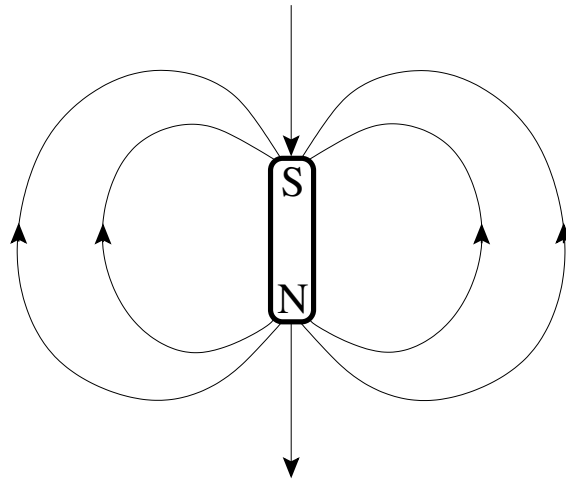


Figure 2.21: Magnetic field-lines generated by a bar magnet.

2.2.9 Magnetic Monopoles

Equation (2.251) immediately suggests that

$$\nabla \cdot \mathbf{B} = 0, \quad (2.263)$$

because the divergence of a curl is identically zero. (See Section A.22.) In other words, the steady magnetic field generated by a pattern of steady circulating electric currents is divergence free. If we integrate the previous equation over a general volume V , bounded by a surface S , making use of the divergence theorem (see Section A.20), then we obtain

$$\oint_S \mathbf{B} \cdot d\mathbf{S} = 0. \quad (2.264)$$

We conclude that the flux of the magnetic field generated by a steady current pattern out of any closed surface is zero. This implies that the magnetic field-lines generated by a steady current pattern are *solenoidal* (see Section A.20) and, consequently, never begin or end.

What about magnetic fields generated by permanent magnets (the modern equivalent of load-stones)? Do they also never begin or end? We know that a conventional bar magnet has both a north and south magnetic pole (like the Earth). If we track the magnetic field-lines with a small compass then they all emanate from the north pole, spread out, and eventually re-converge on the south pole. See Figure 2.21. It appears likely (but we cannot prove it with a compass) that the field-lines inside the magnet connect from the south to the north pole so as to form closed loops that never begin or end.

Can we produce an isolated north or south magnetic pole; for instance, by snapping a bar magnet in two? A compass needle would always point toward an isolated south pole, so this would act like a negative magnetic charge. Likewise, a compass needle would always point away from an isolated north pole, so this would act like a positive magnetic charge. It is clear, from Figure 2.22, that if we take a closed surface S containing an isolated magnetic pole, which is usually termed a

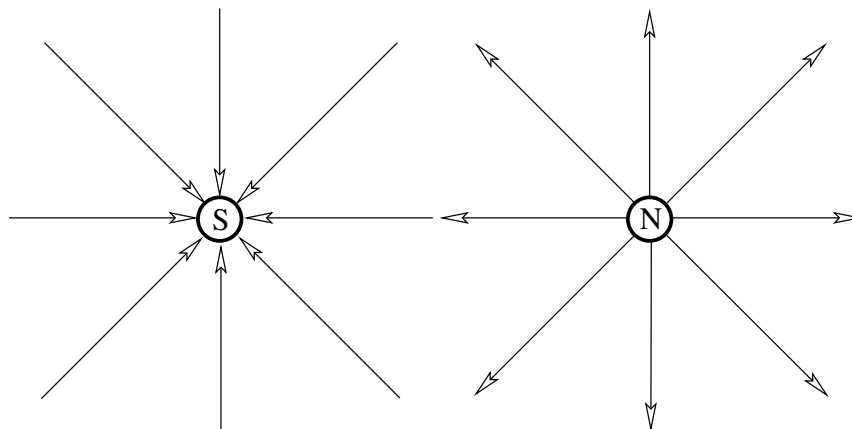


Figure 2.22: Magnetic field-lines generated by magnetic monopoles.

magnetic monopole, then $\oint_S \mathbf{B} \cdot d\mathbf{S} \neq 0$. In fact, the flux will be positive for an isolated north pole, and negative for an isolated south pole. It follows from the divergence theorem (see Section A.20) that if $\oint_S \mathbf{B} \cdot d\mathbf{S} \neq 0$ then $\nabla \cdot \mathbf{B} \neq 0$. Thus, the statement that $\nabla \cdot \mathbf{B} = 0$ is equivalent to the statement that magnetic monopoles do not exist. It is actually quite possible to formulate electromagnetism so as to allow for magnetic monopoles. However, as far as we are aware, there are no magnetic monopoles in the universe. We know that if we try to make a magnetic monopole by snapping a bar magnet in two then we just end up with two smaller bar magnets. If we snap one of these smaller magnets in two then we end up with two even smaller bar magnets. We can continue this process down to the atomic level without ever producing a magnetic monopole. In fact, permanent magnetism is generated by electric currents circulating on the atomic scale, and so this type of magnetism is not fundamentally different to the magnetism generated by macroscopic currents.

In conclusion, all steady magnetic fields in the universe are generated by circulating electric currents of some description. Such fields are solenoidal; that is, they have field-lines that never begin or end, and also satisfy the field equation

$$\nabla \cdot \mathbf{B} = 0. \quad (2.265)$$

We have only proved that $\nabla \cdot \mathbf{B} = 0$ for steady magnetic fields, but, in fact, it turns out that this is also the case for time-dependent fields.

2.2.10 Ampère's Circuital Law

According to Equation (2.251),

$$\nabla \times \mathbf{B} = \nabla \times (\nabla \times \mathbf{A}) \equiv \nabla(\nabla \cdot \mathbf{A}) - \nabla^2 \mathbf{A}, \quad (2.266)$$

where use has been made of Equation (A.187). However, Equation (2.262) indicates that $\nabla \cdot \mathbf{A} = 0$ for a steady current distribution. Hence, the previous equation simplifies to give

$$\nabla \times \mathbf{B} = -\nabla^2 \mathbf{A}. \quad (2.267)$$

The previous equation can be combined with Equation (2.252) to give

$$\nabla \times \mathbf{B}(\mathbf{r}) = -\frac{\mu_0}{4\pi} \int \mathbf{j}(\mathbf{r}') \nabla^2 \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right) dV'. \quad (2.268)$$

However, according to Equation (2.56),

$$\nabla^2 \left(\frac{1}{|\mathbf{r} - \mathbf{r}'|} \right) = -4\pi \delta(\mathbf{r} - \mathbf{r}'), \quad (2.269)$$

so we obtain

$$\nabla \times \mathbf{B}(\mathbf{r}) = \mu_0 \int \mathbf{j}(\mathbf{r}') \delta(\mathbf{r} - \mathbf{r}') dV' = \mu_0 \mathbf{j}(\mathbf{r}), \quad (2.270)$$

where use has been made of Equation (2.47).

The previous equation can be written

$$\nabla \times \mathbf{B} = \mu_0 \mathbf{j}. \quad (2.271)$$

Let us calculate the flux of $\mu_0 \mathbf{j}$ through some surface S , bounded by a loop C . Making use of the previous field equation, as well as the curl theorem (see Section A.22), we obtain

$$\oint_C \mathbf{B} \cdot d\mathbf{r} = \mu_0 \int_S \mathbf{j} \cdot d\mathbf{S}. \quad (2.272)$$

In other words, the line integral of the magnetic field around some loop C is equal to μ_0 multiplied by the net electric current flowing across some surface, S , attached to the loop. This result is known as *Ampère's circuital law*. Note that because the current density associated with a steady current pattern is divergence free [see Equation (2.261)], the net current flowing across any surface attached to C is the same. (See Section A.20.) Of course, when performing the line integral we have to choose an arbitrary sense of circulation around the loop. Once we have done this, any currents that the loop circles in an counter-clockwise direction (looking against the direction of the current) count as positive currents, whereas any currents that the loop circles in a clockwise direction (looking against the direction of the current) count as negative currents.

Let us apply Ampère's circuital law to the trivial case of a circular loop of radius r that lies in the plane perpendicular to a long straight wire, carrying a current I , that passes through its center. By symmetry, we expect the magnetic field to be of the form $\mathbf{B} = B_\theta(r) \mathbf{e}_\theta$, where r, θ, z are right-handed cylindrical polar coordinates defined such that the wire runs along the z -axis. (See Section A.23.) If the chosen sense of circulation around the loop is in the direction of increasing θ then I counts as a positive current. Thus, Equation (2.272) yields

$$2\pi r B_\theta(r) = \mu_0 I, \quad (2.273)$$

or

$$B_\theta = \frac{\mu_0 I}{2\pi r}, \quad (2.274)$$

which is equivalent to Equation (2.206).

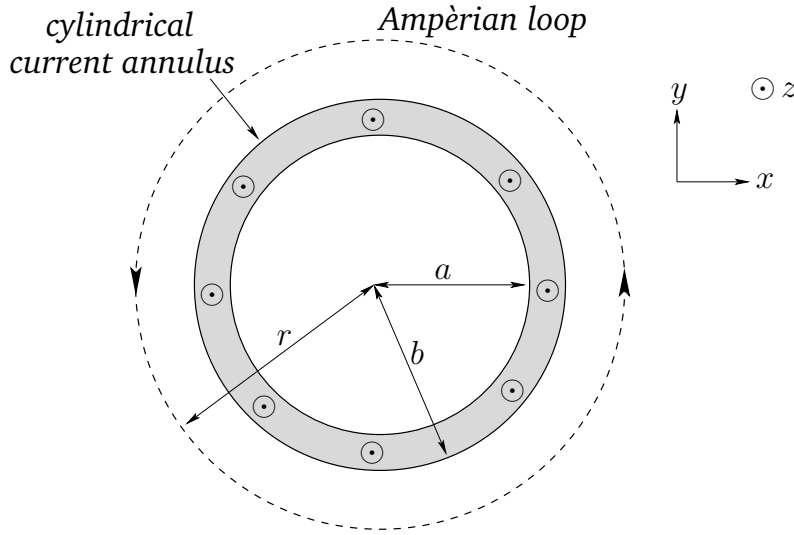


Figure 2.23: An example use of Ampère's circuital law.

As another example of the use of Ampère's circuital law, let us calculate the magnetic field generated by a cylindrical current annulus of inner radius a , and outer radius b , co-axial with the z -axis, and carrying a uniformly distributed z -directed current I . By symmetry, and also by analogy with the magnetic field generated by a straight wire, we expect the current distribution to generate a magnetic field of the form $\mathbf{B} = B_\theta(r) \mathbf{e}_\theta$, where r, θ, z are right-handed cylindrical polar coordinates. (See Section A.23.) Let us apply Ampère's circuital law to an imaginary circular loop in the $x-y$ plane, of radius r , centered on the z -axis. See Figure 2.23. Such a loop is generally known as an *Ampèrian loop*. As before, if the chosen sense of circulation around the loop is in the direction of increasing θ then I counts as a positive current. According to Ampère's circuital law, the line integral of the magnetic field around the loop is equal to the current passing through the plane of the loop, multiplied by μ_0 . The line integral is easy to calculate because the magnetic field is everywhere tangential to the loop. We obtain

$$2\pi r B_\theta(r) = \mu_0 I(r),$$

where $I(r)$ is the current that passes through an Ampèrian loop of radius r . Simple arguments involving proportion reveal that

$$I(r) = \begin{cases} 0 & r < a \\ [(r^2 - a^2)/(b^2 - a^2)] I & a \leq r \leq b \\ I & b < r \end{cases} . \quad (2.275)$$

Hence,

$$B_\theta(r) = \begin{cases} 0 & r < a \\ [\mu_0 I/(2\pi r)] [(r^2 - a^2)/(b^2 - a^2)] & a \leq r \leq b \\ \mu_0 I/(2\pi r) & b < r \end{cases} . \quad (2.276)$$

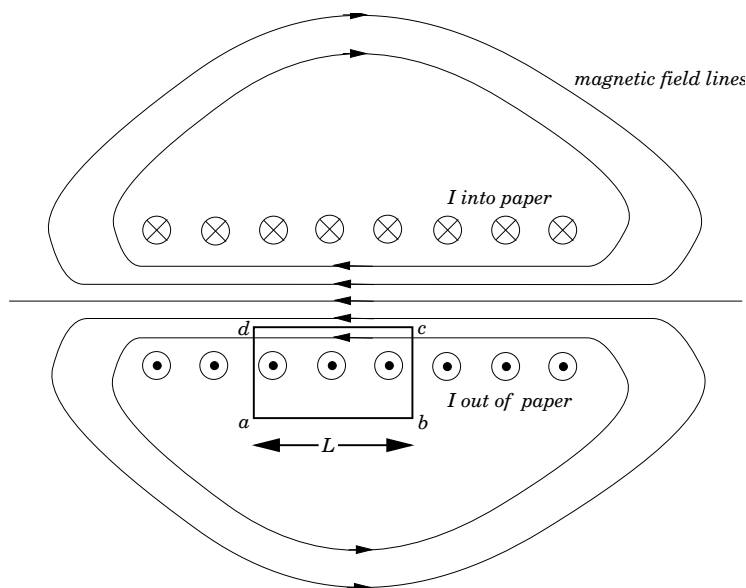


Figure 2.24: A solenoid.

2.2.11 Magnetic Field of a Solenoid

A *solenoid* is a tightly wound, helical coil of wire whose diameter is small compared to its length. The magnetic field generated in the center, or core, of a current-carrying solenoid is essentially uniform, and is directed along the axis of the solenoid. Outside the solenoid, the magnetic field is far weaker. Figure 2.24 shows (rather schematically) the magnetic field generated by a typical solenoid. The solenoid is wound from a single helical wire that carries a current I . The winding is sufficiently tight that each turn of the solenoid is well approximated as a circular wire loop, lying in the plane perpendicular to the axis of the solenoid, that carries a current I . Suppose that there are N such turns per unit axial length of the solenoid. What is the magnitude of the magnetic field in the core of the solenoid?

In order to answer this question, let us apply Ampère's circuital law to the rectangular loop $abcd$. We must first find the line integral of the magnetic field around $abcd$. Along bc and da the magnetic field is essentially perpendicular to the loop, so there is no contribution to the line integral from these sections of the loop. Along cd the magnetic field is approximately uniform, of magnitude B , say, and is directed parallel to the loop. Thus, the contribution to the line integral from this section of the loop is BL , where L is the length of cd . Along ab the magnetic field-strength is essentially negligible, so this section of the loop makes no contribution to the line integral. It follows that the line integral of the magnetic field around $abcd$ is simply

$$w = BL. \quad (2.277)$$

By Ampère's circuital law, this line integral is equal to μ_0 multiplied by the algebraic sum of the currents that pass through the plane of the loop $abcd$. Because the length of the loop along the axis of the solenoid is L , the loop encloses NL turns of the solenoid, each of which carries a current I . Thus, the total current that passes through the plane of the loop is NLI . This current counts

as a positive current, because if we look against the direction of the currents flowing in each turn (i.e., into the page in the figure) then the loop $abcd$ circulates these currents in a counter-clockwise direction. Ampère's circuital law yields

$$BL = \mu_0 NLI, \quad (2.278)$$

which reduces to

$$B = \mu_0 NI. \quad (2.279)$$

Thus, the magnetic field in the core of a solenoid is directly proportional to the product of the current flowing around the solenoid and the number of turns per unit length of the solenoid. This result is exact in the limit in which the length of the solenoid is very much greater than its diameter.

2.3 Magnetic Induction

2.3.1 Faraday's Law

The phenomenon of magnetic induction plays a crucial role in three very useful electrical devices; the electric generator (see Section 2.3.10), the electric motor (see Section 2.3.12), and the transformer (see Section 2.3.13). Without these devices, modern life would be impossible in its present form. Magnetic induction was discovered in 1830 by Michael Faraday. Joseph Henry independently made the same discovery at about the same time. Both physicists were intrigued by the fact that an electric current flowing around a circuit can generate a magnetic field. Surely, they reasoned, if an electric current can generate a magnetic field then a magnetic field must somehow be able to generate an electric current. However, it took many years of fruitless experimentation before they were able to find the essential ingredient that allows a magnetic field to generate an electric current. This ingredient is time variation.

Prior to 1830, the only known way in which to cause an electric current to flow through a conducting wire was to connect the ends of the wire to the positive and negative terminals of a battery. We measure a battery's ability to push current down a wire in terms of its voltage, by which we mean the voltage difference between its positive and negative terminals. Of course, volts are the units used to measure electric scalar potential, so when we talk about a 6V battery, what we are really saying is that the difference in electric scalar potential between its positive and negative terminals is six volts. This insight allows us to write

$$V = \phi(\oplus) - \phi(\ominus) = - \int_{\oplus}^{\ominus} \nabla\phi \cdot d\mathbf{r} = \int_{\oplus}^{\ominus} \mathbf{E} \cdot d\mathbf{r}, \quad (2.280)$$

where V is the battery voltage, \oplus denotes the positive terminal, \ominus the negative terminal, and $d\mathbf{r}$ is an element of length along the wire. Of course, the previous equation is a direct consequence of $\mathbf{E} = -\nabla\phi$. [See Equation (2.17) and Section A.18.] Clearly, a voltage difference between two ends of a wire attached to a battery implies the presence of a longitudinal electric field that pushes electric charges along the wire. This field is directed from the positive terminal of the battery to the negative terminal, and is, therefore, such as to force electrons to flow through the wire from the

negative to the positive terminal. As expected, this implies that a net positive current flows from the positive to the negative terminal. The fact that \mathbf{E} is a conservative field (i.e., $\mathbf{E} = -\nabla\phi$) ensures that the voltage difference, V , is independent of the path of the wire between the terminals. In other words, two different wires attached to the same battery develop identical voltage differences.

Let us now consider a closed loop of wire (with no battery). The voltage around such a loop, which is sometimes called the *electromotive force*, or *emf*, is

$$V = \oint \mathbf{E} \cdot d\mathbf{r} = 0. \quad (2.281)$$

The fact that the right-hand side of the previous equation is zero is a direct consequence of the field equation $\nabla \times \mathbf{E} = -\nabla \times \nabla\phi = \mathbf{0}$ and the curl theorem. [See Equations (2.17) and (2.25), and Section A.22.] We conclude that, because \mathbf{E} is a conservative field (i.e., $\mathbf{E} = -\nabla\phi$), the emf around a closed loop of wire is automatically zero, and so there is no current flow around such a loop.

However, in 1830, Michael Faraday discovered that a changing magnetic field can cause a current to flow around a closed loop of wire (in the absence of a battery). Of course, if current flows around the loop then there must be an emf. In other words,

$$V = \oint \mathbf{E} \cdot d\mathbf{r} \neq 0, \quad (2.282)$$

which immediately implies that \mathbf{E} is not a conservative field, and that $\nabla \times \mathbf{E} \neq \mathbf{0}$. Clearly, we are going to have to modify some of our ideas regarding electric fields.

Faraday continued his experiments, and found that another way of generating an emf around a loop of wire is to keep the magnetic field constant and to move the loop. (See Section 2.3.9.) Eventually, Faraday was able to formulate a law that accounted for all of his experiments; the emf generated around a loop of wire in a magnetic field is proportional to the rate of change of the flux of the magnetic field through the loop. (See Section A.20.) Thus, if the loop is denoted C , and S is some surface attached to the loop, then Faraday's experiments can be summed up by writing

$$V = \oint_C \mathbf{E} \cdot d\mathbf{r} = A \frac{\partial}{\partial t} \int_S \mathbf{B} \cdot d\mathbf{S}, \quad (2.283)$$

where A is a constant of proportionality. Thus, the changing flux of the magnetic field passing through the loop generates an electric field directed around the loop. This process is known as *magnetic induction*.

SI units have been carefully chosen so as to make $|A| = 1$ in the previous equation. So, the only question that we now have to answer is whether $A = +1$ or $A = -1$. In other words, we need to decide which way around the loop the induced emf drives the current. We possess a general principle, known as *Le Chatelier's principle*, that allows us to answer such questions. According to Le Chatelier's principle, every change in a physical system generates a reaction that acts to minimize the change. Essentially, this implies that the universe is stable to small perturbations. When Le Chatelier's principle is applied to the particular case of magnetic induction, it is usually called *Lenz's law*, after Emil Lenz who formulated it in 1834. According to Lenz's law, the current induced by an emf around a closed loop is always such that the magnetic field it produces acts to

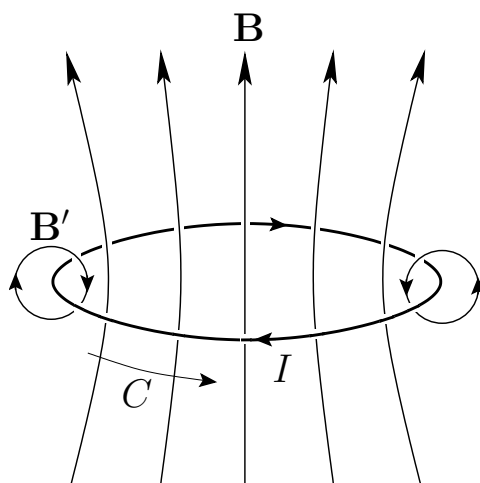


Figure 2.25: Lenz's law.

counteract the change in magnetic flux that generates the emf. From Figure 2.25, it is clear that if the magnetic field \mathbf{B} is increasing and the current I circulates clockwise (as seen from above) then the current generates a field \mathbf{B}' that opposes the increase in the magnetic flux through the loop, in accordance with Lenz's law. The direction of the current is opposite to the sense of circulation of the current loop C , as determined by the right-hand rule (assuming that the flux of \mathbf{B} through the loop is positive), so this implies that $A = -1$ in Equation (2.283). Thus, Faraday's law takes the form

$$V = \oint_C \mathbf{E} \cdot d\mathbf{r} = -\frac{\partial}{\partial t} \int_S \mathbf{B} \cdot d\mathbf{S} = -\frac{d\Phi}{dt}, \quad (2.284)$$

where $\Phi = \int_S \mathbf{B} \cdot d\mathbf{S}$ is the magnetic flux through the loop.

Experimentally, Faraday's law is found to correctly predict the emf (i.e., $\oint \mathbf{E} \cdot d\mathbf{r}$) generated around any wire loop, irrespective of the position or shape of the loop. It is reasonable to assume that the same emf would be generated in the absence of the wire (of course, no current would flow in this case). We conclude that Equation (2.284) is valid for any closed loop C . Now, if Faraday's law is to make sense then it must hold for all surfaces, S , attached to the loop, C . Clearly, if the flux of the magnetic field through the loop depends on the surface upon which it is evaluated then Faraday's law is going to predict different emfs for different surfaces. Because there is no preferred surface for a general non-coplanar loop, this would not make any sense. The condition for the flux of the magnetic field, $\int_S \mathbf{B} \cdot d\mathbf{S}$, to depend only on the loop C to which the surface S is attached, and not on the nature of the surface itself, is

$$\oint_{S'} \mathbf{B} \cdot d\mathbf{S}' = 0, \quad (2.285)$$

for any closed surface S' . (See Section A.20.)

Faraday's law, Equation (2.284), can be converted into a field equation using the curl theorem. (See Section A.22.) We obtain

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}. \quad (2.286)$$

This field equation describes how a changing magnetic field generates an electric field. The divergence theorem (see Section A.20) applied to Equation (2.285) gives the familiar field equation

$$\nabla \cdot \mathbf{B} = 0. \quad (2.287)$$

[See Equation (2.263).] This equation ensures that the magnetic flux through a loop is a well defined quantity.

The divergence of Equation (2.286) yields

$$\frac{\partial(\nabla \cdot \mathbf{B})}{\partial t} = 0, \quad (2.288)$$

because $\nabla \cdot \nabla \times \mathbf{E} \equiv 0$. (See Section A.22.) Thus, the field equation (2.286) actually demands that the divergence of the magnetic field must be constant in time for self-consistency (this implies that the flux of the magnetic field through a loop need not be a well defined quantity, as long as its time derivative is well defined). However, a constant non-solenoidal magnetic field can only be generated by magnetic monopoles, and magnetic monopoles do not exist (as far as we are aware). (See Section 2.2.9.) Hence, $\nabla \cdot \mathbf{B} = 0$.

As an example of the use of Faraday's law, let us calculate the electric field generated by a decaying magnetic field of the form $\mathbf{B} = B_z(r, t) \mathbf{e}_z$, where

$$B_z(r, t) = \begin{cases} B_0 \exp(-t/\tau) & r \leq a \\ 0 & r > a \end{cases}, \quad (2.289)$$

and r is a cylindrical polar coordinate. (See Section A.23.) Here, B_0 and τ are positive constants. By symmetry, we expect an induced electric field of the form $\mathbf{E}(r, t)$. We also expect $\nabla \cdot \mathbf{E} = 0$, because there are no electric charges in the problem. [See Equation (2.54).] This rules out a radial electric field. We can also rule out a z -directed electric field, because $\nabla \times [E_z(r) \mathbf{e}_z] = -(\partial E_z / \partial r) \mathbf{e}_\theta$, and we require $\nabla \times \mathbf{E} \propto \mathbf{B} \propto \mathbf{e}_z$. Hence, the induced electric field must be of the form $\mathbf{E}(r, t) = E_\theta(r, t) \mathbf{e}_\theta$. Now, according to Faraday's law, (2.284), the line integral of the electric field around some closed loop is equal to minus the rate of change of the magnetic flux passing through the loop. If we choose a loop that is a circle of radius r in the x - y plane then we have

$$2\pi r E_\theta(r, t) = -\frac{d\Phi}{dt}, \quad (2.290)$$

where Φ is the flux of the magnetic field (in the $+z$ direction) passing through a circular loop of radius r . It is evident that

$$\Phi(r, t) = \begin{cases} \pi r^2 B_0 \exp(-t/\tau) & r \leq a \\ \pi a^2 B_0 \exp(-t/\tau) & r > a \end{cases}. \quad (2.291)$$

Hence,

$$E_\theta(r, t) = \begin{cases} (B_0/2\tau) r \exp(-t/\tau) & r \leq a \\ (B_0/2\tau) (a^2/r) \exp(-t/\tau) & r > a \end{cases}. \quad (2.292)$$

2.3.2 Electric Scalar Potential

We now have a problem. We can only write the electric field in terms of a scalar potential (i.e., $\mathbf{E} = -\nabla\phi$) provided that $\nabla \times \mathbf{E} = \mathbf{0}$. This follows because $\nabla \times \nabla\phi \equiv \mathbf{0}$. (See Section A.22.) However, we have just discovered that the curl of the electric field is non-zero in the presence of a changing magnetic field. In other words, \mathbf{E} is not, in general, a conservative field. Does this mean that we have to abandon the concept of electric scalar potential? Fortunately, it does not. It is still possible to define a scalar potential that is physically meaningful.

Let us start from the field equation

$$\nabla \cdot \mathbf{B} = 0, \quad (2.293)$$

which is valid for both time-varying and constant magnetic fields. Because the magnetic field is solenoidal, we can write it as the curl of a vector potential:

$$\mathbf{B} = \nabla \times \mathbf{A}. \quad (2.294)$$

[See Equation (2.251)]. This follows because $\nabla \cdot (\nabla \times \mathbf{A}) \equiv 0$. (See Section A.22.) So, there is no problem with the vector potential in the presence of time-varying fields. Let us substitute Equation (2.294) into the field equation (2.286). We obtain

$$\nabla \times \mathbf{E} = -\frac{\partial (\nabla \times \mathbf{A})}{\partial t}, \quad (2.295)$$

which can be written

$$\nabla \times \left(\mathbf{E} + \frac{\partial \mathbf{A}}{\partial t} \right) = \mathbf{0}. \quad (2.296)$$

Now, we know that a curl-free vector field can always be expressed as the gradient of a scalar potential (see Section A.22), so let us write

$$\mathbf{E} + \frac{\partial \mathbf{A}}{\partial t} = -\nabla\phi, \quad (2.297)$$

or

$$\mathbf{E} = -\nabla\phi - \frac{\partial \mathbf{A}}{\partial t}. \quad (2.298)$$

This equation implies that the electric scalar potential, ϕ , only describes the conservative electric field generated by electric charges. The electric field induced by time-varying magnetic fields is non-conservative, and is described by the magnetic vector potential, \mathbf{A} .

2.3.3 Gauge Invariance

As we saw in the previous section, electric and magnetic fields can be written in terms of scalar and vector potentials, as follows:

$$\mathbf{E} = -\nabla\phi - \frac{\partial \mathbf{A}}{\partial t}, \quad (2.299)$$

$$\mathbf{B} = \nabla \times \mathbf{A}. \quad (2.300)$$

However, this prescription is not unique. There are many different potentials that can generate the same fields. This phenomenon is known as *gauge invariance*. The most general transformation that leaves the \mathbf{E} and \mathbf{B} fields unchanged in Equations (2.299) and (2.300) is

$$\phi \rightarrow \phi + \frac{\partial \psi}{\partial t}, \quad (2.301)$$

$$\mathbf{A} \rightarrow \mathbf{A} - \nabla \psi, \quad (2.302)$$

where $\psi(\mathbf{r}, t)$ is a general scalar field known as the *gauge field*. A particular choice of the gauge field is termed a choice of the gauge.

We are free to choose the gauge so as to make our equations as simple as possible. As before, the most sensible gauge for the scalar potential is to set it to zero at infinity:

$$\phi(\mathbf{r}, t) \rightarrow 0 \quad \text{as } |\mathbf{r}| \rightarrow \infty. \quad (2.303)$$

For steady fields, we found that

$$\nabla \cdot \mathbf{A} = 0. \quad (2.304)$$

[See Equation (2.262).] This choice is known as the *Coulomb gauge*. We can still use this gauge for time-varying fields.

Equation (2.299) can be combined with the field equation [see Equation (2.54)]

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0} \quad (2.305)$$

(which remains valid for time-varying fields) to give

$$-\nabla^2 \phi - \frac{\partial (\nabla \cdot \mathbf{A})}{\partial t} = \frac{\rho}{\epsilon_0}. \quad (2.306)$$

(See Section A.21.) With the Coulomb gauge, $\nabla \cdot \mathbf{A} = 0$, the previous expression reduces to

$$\nabla^2 \phi = -\frac{\rho}{\epsilon_0}, \quad (2.307)$$

which is just Poisson's equation. [See Equation (2.99).] Thus, we can immediately write down an expression for the scalar potential generated by time-varying fields. It is exactly analogous to our previous expression for the scalar potential generated by steady fields:

$$\phi(\mathbf{r}, t) = \frac{1}{4\pi \epsilon_0} \int \frac{\rho(\mathbf{r}', t)}{|\mathbf{r} - \mathbf{r}'|} dV'. \quad (2.308)$$

[See Equation (2.18).] However, this apparently simple result is extremely deceptive. Equation (2.308) is a typical action at a distance law. If the charge density changes suddenly at \mathbf{r}' then the potential at \mathbf{r} responds immediately. However, special relativity forbids information from propagating faster than the speed of light in vacuum, because this would violate causality. (See Section 3.2.10.)

How can these two statements be reconciled? The crucial point is that the scalar potential cannot be measured directly, it can only be inferred from the electric field. In the time dependent case, there are two parts to the electric field; that part that comes from the scalar potential, and that part that comes from the vector potential. [See Equation (2.299).] So, if the scalar potential in some region responds immediately to some distance rearrangement of charge density then it does not necessarily follow that the electric field also has an immediate response. What actually happens is that the change in the part of the electric field that comes from the scalar potential is balanced by an equal and opposite change in the part that comes from the vector potential, so that the overall electric field remains unchanged. This state of affairs persists at least until sufficient time has elapsed for a light signal to travel from the distant charges to the region in question. Thus, causality is not violated, because it is the electric field, and not the scalar potential, that carries physically accessible information.

It is clear that the apparent action at a distance nature of Equation (2.308) is highly misleading. This suggests, very strongly, that the Coulomb gauge is not the optimum gauge in the time dependent case. A more sensible choice is the so-called *Lorenz gauge*:

$$\nabla \cdot \mathbf{A} = -\epsilon_0 \mu_0 \frac{\partial \phi}{\partial t}. \quad (2.309)$$

Substituting the Lorenz gauge into Equation (2.306), we obtain

$$\epsilon_0 \mu_0 \frac{\partial^2 \phi}{\partial t^2} - \nabla^2 \phi = \frac{\rho}{\epsilon_0}. \quad (2.310)$$

It turns out that this is a three-dimensional wave equation in which information propagates at the speed of light in vacuum. (See Section 2.4.4.) Thus, the Lorenz gauge makes manifest the fact that information carried by electric and magnetic fields propagates at the velocity of light in vacuum, which implies that causality is not violated.

2.3.4 Inductance

We have already learned about the concepts of voltage, resistance, and capacitance. Let us now investigate the concept of *inductance*. Electrical engineers like to reduce all pieces of electrical circuitry to an *equivalent circuit* consisting of pure voltage sources, pure inductors, pure capacitors, and pure resistors. Hence, once we understand inductors, we shall be ready to apply the laws of electromagnetism to general electrical circuits.

Consider two stationary loops of wire, labeled 1 and 2. See Figure 2.26. Let us run a steady current I_1 around the first loop to produce a magnetic field \mathbf{B}_1 . Some of the field-lines of \mathbf{B}_1 will pass through the second loop. Let Φ_2 be the flux of \mathbf{B}_1 through loop 2,

$$\Phi_2 = \int_{\text{loop 2}} \mathbf{B}_1 \cdot d\mathbf{S}_2, \quad (2.311)$$

where $d\mathbf{S}_2$ is a surface element of loop 2. This flux is generally quite difficult to calculate exactly (unless the two loops have a particularly simple geometry). However, we can infer from the Biot-

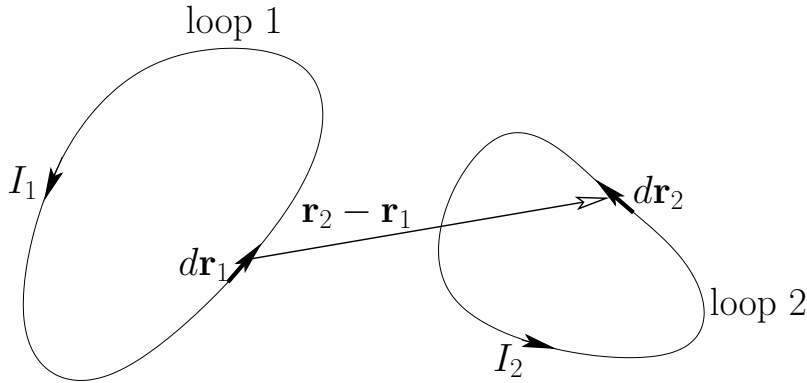


Figure 2.26: Two current-carrying loops.

Savart law [see Equation (2.226)],

$$\mathbf{B}_1(\mathbf{r}) = \frac{\mu_0 I_1}{4\pi} \oint_{\text{loop 1}} \frac{d\mathbf{r}_1 \times (\mathbf{r} - \mathbf{r}_1)}{|\mathbf{r} - \mathbf{r}_1|^3}, \quad (2.312)$$

that the magnitude of \mathbf{B}_1 is proportional to the current I_1 . Here, $d\mathbf{r}_1$ is a line element of loop 1 located at displacement \mathbf{r}_1 . It follows that the flux Φ_2 must also be proportional to I_1 . Thus, we can write

$$\Phi_2 = M_{21} I_1, \quad (2.313)$$

where M_{21} is a constant of proportionality. This constant is termed the *mutual inductance* of the two loops.

Let us write the magnetic field \mathbf{B}_1 in terms of a vector potential \mathbf{A}_1 , so that

$$\mathbf{B}_1 = \nabla \times \mathbf{A}_1. \quad (2.314)$$

It follows from the curl theorem (see Section A.22) that

$$\Phi_2 = \int_{\text{loop 2}} \mathbf{B}_1 \cdot d\mathbf{S}_2 = \int_{\text{loop 2}} \nabla \times \mathbf{A}_1 \cdot d\mathbf{S}_2 = \oint_{\text{loop 2}} \mathbf{A}_1 \cdot d\mathbf{r}_2, \quad (2.315)$$

where $d\mathbf{r}_2$ is a line element of loop 2. However, we know that

$$\mathbf{A}_1(\mathbf{r}) = \frac{\mu_0 I_1}{4\pi} \oint_{\text{loop 1}} \frac{d\mathbf{r}_1}{|\mathbf{r} - \mathbf{r}_1|}. \quad (2.316)$$

The previous equation is just a special case of the more general result [see Equation (2.252)],

$$\mathbf{A}_1(\mathbf{r}) = \frac{\mu_0}{4\pi} \int \frac{\mathbf{j}(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} dV', \quad (2.317)$$

for $\mathbf{j}(\mathbf{r}_1) = d\mathbf{r}_1 I_1 / (dr_1 dA)$ and $dV' = dr_1 dA$, where dA is the cross-sectional area of loop 1. Thus,

$$\Phi_2 = \frac{\mu_0 I_1}{4\pi} \oint_{\text{loop 1}} \oint_{\text{loop 2}} \frac{d\mathbf{r}_1 \cdot d\mathbf{r}_2}{|\mathbf{r}_2 - \mathbf{r}_1|}, \quad (2.318)$$

where \mathbf{r}_2 is the position vector of the line element $d\mathbf{r}_2$ of loop 2, which implies that

$$M_{21} = \frac{\mu_0}{4\pi} \oint_{\text{loop 1}} \oint_{\text{loop 2}} \frac{d\mathbf{r}_1 \cdot d\mathbf{r}_2}{|\mathbf{r}_2 - \mathbf{r}_1|}. \quad (2.319)$$

In fact, mutual inductances are rarely worked out using the previous formula, because it is usually far too difficult. However, this formula, which is known as the *Neumann formula*, tells us two important things. Firstly, the mutual inductance of two current loops is a purely geometric quantity, having to do with the sizes, shapes, and relative orientations of the loops. Secondly, the integral is unchanged if we switch the roles of loops 1 and 2. In other words,

$$M_{21} = M_{12}. \quad (2.320)$$

Hence, we can drop the subscripts, and just call both of these quantities M . This result implies that no matter what the shapes and relative positions of the two loops, the magnetic flux through loop 2 when a current I runs around loop 1 is exactly the same as the flux through loop 1 when the same current runs around loop 2.

We have seen that a current I flowing around some wire loop, 1, generates a magnetic flux linking some other loop, 2. However, flux is also generated through the first loop. As before, the magnetic field, and, therefore, the flux, Φ , is proportional to the current, so we can write

$$\Phi = LI. \quad (2.321)$$

The constant of proportionality L is called the *self inductance*. Like M it only depends on the geometry of the loop.

The SI unit of inductance is the *henry* (H), which is equivalent to a volt-second per ampere. The henry, like the farad, is a rather unwieldy unit, because inductors in electrical circuits typically have a inductances of order a micro-henry.

2.3.5 Self Inductance

Consider a long, uniformly wound, cylindrical solenoid of length l , and radius r , that has N turns per unit length, and carries a current I . The longitudinal (i.e., directed along the axis of the solenoid) magnetic field within the solenoid is approximately uniform, and is given by

$$B = \mu_0 N I. \quad (2.322)$$

(See Section 2.2.11.) The magnetic flux passing through each turn of the solenoid wire is $B\pi r^2 = \mu_0 N I \pi r^2$. Thus, the total flux passing through the solenoid wire, which has Nl turns, is

$$\Phi = Nl\mu_0 N I \pi r^2. \quad (2.323)$$

Hence, the self inductance of the solenoid is

$$L = \frac{\Phi}{I} = \mu_0 N^2 \pi r^2 l. \quad (2.324)$$

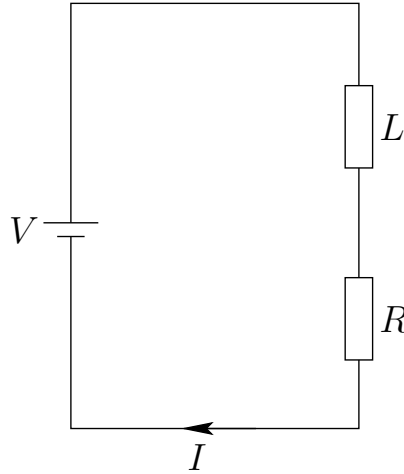


Figure 2.27: The equivalent circuit of a solenoid connected to a battery.

Note that the self inductance only depends on geometric quantities, such as the number of turns per unit length of the solenoid, and the cross-sectional area of the turns.

Suppose that the current I flowing through the solenoid changes. A change in the current implies a change in the magnetic flux linking the solenoid wire, because $\Phi = LI$. According to Faraday's law, this change generates an emf in the wire. By Lenz's law, the emf is such as to oppose the change in the current; that is, it is a *back-emf*. Thus, we can write

$$\mathcal{E} = -\frac{d\Phi}{dt} = -L \frac{dI}{dt}, \quad (2.325)$$

where \mathcal{E} is the generated back-emf. [See Equation (2.284).]

Suppose that our solenoid has an electrical resistance R . Let us connect the ends of the solenoid across the terminals of a battery of constant voltage V . The equivalent circuit is shown in Figure 2.27. The inductance and resistance of the solenoid are represented by a perfect inductor, L , and a perfect resistor, R , connected in series. The voltage drop across the inductor and resistor is equal to the voltage of the battery, V . The voltage drop across the resistor is simply IR (see Section 2.1.11), whereas the voltage drop across the inductor (i.e., minus the back-emf) is $L dI/dt$. Here, I is the current flowing through the solenoid. It follows that

$$V = IR + L \frac{dI}{dt}. \quad (2.326)$$

This is a differential equation for the current I . We can rearrange it to give

$$\frac{dI}{dt} + \frac{R}{L} I = \frac{V}{L}. \quad (2.327)$$

The general solution is

$$I(t) = \frac{V}{R} + k \exp\left(-\frac{Rt}{L}\right). \quad (2.328)$$

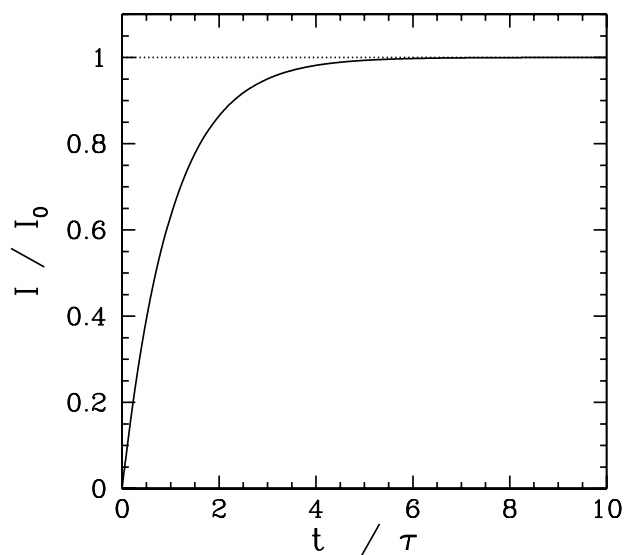


Figure 2.28: Typical current rise profile in a circuit of the type shown in Figure 2.27. Here, $I_0 = V/R$ and $\tau = L/R$.

The constant k is fixed by the initial conditions. Suppose that the battery is connected at time $t = 0$, when $I = 0$. It follows that $k = -V/R$, so that

$$I(t) = \frac{V}{R} \left[1 - \exp\left(-\frac{Rt}{L}\right) \right]. \quad (2.329)$$

This curve is shown in Figure 2.28. It can be seen that, after the battery is connected, the current ramps up, and attains its steady-state value V/R (which comes from Ohm's law), on the characteristic timescale

$$\tau = \frac{L}{R}. \quad (2.330)$$

To be more exact, the current has risen to approximately 63% of its final value at time $t = \tau$, and to more than 99% of its final value at time $t = 5\tau$. The timescale τ is sometimes called the *time constant* of the circuit, or (somewhat unimaginatively) the *L over R time* of the circuit. We conclude that it takes a finite time to establish a steady current flowing through a solenoid.

2.3.6 RC Circuits

Let us now discuss a topic that, admittedly, has nothing whatsoever to do with inductors, but is mathematically so similar to the topic just discussed that it seems sensible to consider it at this point.

Consider a circuit in which a battery of emf V is connected in series with a capacitor of capacitance C , and a resistor of resistance R . For fairly obvious reasons, such a circuit is generally referred to as an *RC circuit*. In a steady state, the charge on the positive plate of the capacitor is

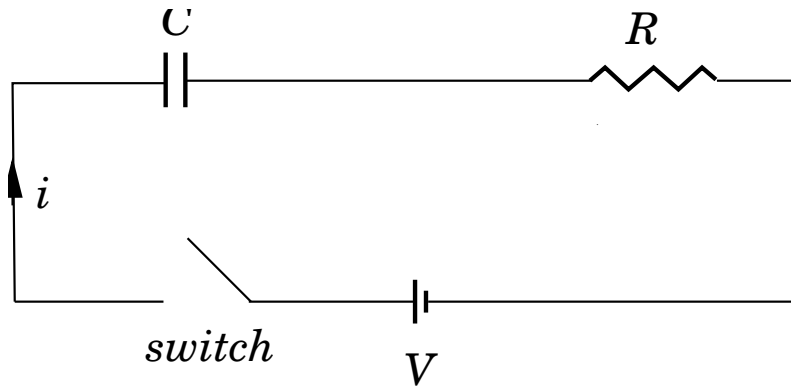


Figure 2.29: An RC circuit with a switch.

given by $Q = CV$, and zero current flows around the circuit (because current cannot flow across the insulating gap between the capacitor plates).

Let us now introduce a switch into the circuit, as shown in Figure 2.29. Suppose that the switch is initially open, but is suddenly closed at $t = 0$. It is assumed that the capacitor plates are uncharged when the switch is thrown. We expect a transient current, i , to flow around the circuit until the charge, q , on the positive plate of the capacitor attains its final steady-state value, $Q = CV$. But, how long does this process take?

The potential difference, v , between the positive and negative plates of the capacitor is given by

$$v = V - iR. \quad (2.331)$$

In other words, the potential difference between the plates is the emf of the battery minus the potential drop across the resistor. The charge, q , on the positive plate of the capacitor is written

$$q = Cv = Q - iRC, \quad (2.332)$$

where $Q = CV$ is the final charge. Now, if i is the instantaneous current flowing around the circuit then, in a short time interval dt , the charge on the positive plate of the capacitor increases by a small amount $dq = i dt$ (because all of the charge that flows around the circuit must accumulate on the plates of the capacitor). It follows that

$$i = \frac{dq}{dt}. \quad (2.333)$$

Thus, the instantaneous current flowing around the circuit is numerically equal to the rate at which the charge accumulated on the positive plate of the capacitor increases with time. Equations (2.332) and (2.333) can be combined together to give

$$\frac{dq'}{dt} = -\frac{q'}{RC}, \quad (2.334)$$

where

$$q' = q - Q. \quad (2.335)$$

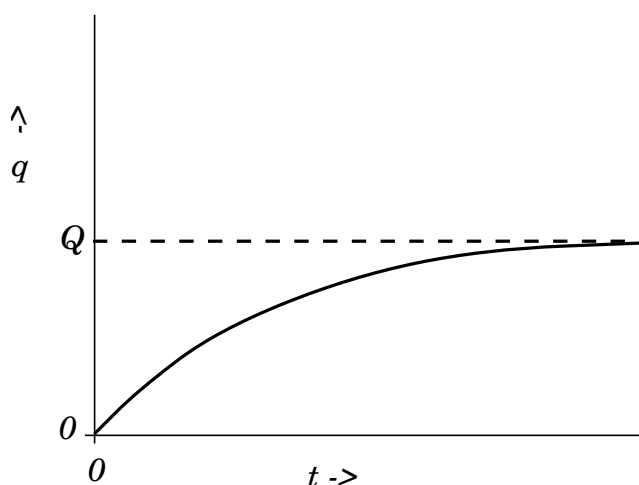


Figure 2.30: Sketch of the charging phase in an RC circuit switched on at $t = 0$.

At $t = 0$, just after the switch is closed, the charge on the positive plate of the capacitor is zero, so

$$q'(t = 0) = -Q. \quad (2.336)$$

Integration of Equation (2.334), subject to the initial condition (2.336), yields

$$q'(t) = -Q \exp\left(-\frac{t}{RC}\right). \quad (2.337)$$

It follows from Equation (2.335) that

$$q(t) = Q \left[1 - \exp\left(-\frac{t}{RC}\right) \right] \quad (2.338)$$

The previous expression specifies the charge, q , on the positive plate of the capacitor a time interval, t , after the switch is closed (at time $t = 0$). The variation of the charge with time is sketched in Figure 2.30. It can be seen that, when the switch is closed, the charge q on the positive plate of the capacitor does not suddenly jump up to its final value, $Q = CV$. Instead, the charge increases smoothly from zero, and gradually asymptotes to its final value. The charge has risen to approximately 63% of its final value a time

$$\tau = RC \quad (2.339)$$

after the switch is closed. By the time $t = 5\tau$, the charge has risen to more than 99% of its final value. Thus, $\tau = RC$ is a good measure of how long after the switch is closed it takes for the capacitor to fully charge up. The quantity τ is termed the time constant, or, somewhat unimaginatively, the RC time, of the circuit.

According to Equations (2.333) and (2.334),

$$i = \frac{dq}{dt} = \frac{dq'}{dt} = -\frac{q'}{RC}. \quad (2.340)$$

It follows from Equation (2.337) that

$$i(t) = I \exp\left(-\frac{t}{RC}\right), \quad (2.341)$$

where $I = V/R$. The previous expression specifies the current, i , flowing around the circuit a time interval, t , after the switch is closed (at time $t = 0$). It can be seen that, immediately after the switch is thrown, the current, $I = V/R$, that flows in the circuit is that which would flow if the capacitor were replaced by a conducting wire. However, this current is only transient, and rapidly decays away to a negligible value. After one RC time, the current has decayed to 37% of its initial value. After five RC times, the current has decayed to less than 1% of its initial value. It is interesting to note that, for a short instant of time, just after the switch is closed, the current in the circuit acts as if there is no insulating gap between the capacitor plates. It essentially takes an RC time for the information about the break in the circuit to propagate around the circuit, and cause the current to stop flowing.

2.3.7 Mutual Inductance

Consider two long, thin, cylindrical solenoids, one wound on top of the other. The common length of each solenoid is l , and the common radius is r . Suppose that the inner solenoid has N_1 turns per unit length, and carries a current I_1 . The magnetic field generated within the inner solenoid is $B_1 = \mu_0 N_1 I_1$. [See Equation (2.322).] The magnetic flux passing through each turn of the outer solenoid is $\mu_0 N_1 I_1 \pi r^2$, and the total flux linking the outer solenoid is therefore $\Phi_2 = N_2 l \mu_0 N_1 I_1 \pi r^2$, where N_2 is the number of turns per unit length of the outer solenoid. It follows that the mutual inductance of the two solenoids, defined $\Phi_2 = M I_1$ [see Equation (2.313)] is given by

$$M = \mu_0 N_1 N_2 \pi r^2 l. \quad (2.342)$$

Recall that the self inductance of the inner solenoid is

$$L_1 = \mu_0 N_1^2 \pi r^2 l, \quad (2.343)$$

and that of the outer solenoid is

$$L_2 = \mu_0 N_2^2 \pi r^2 l. \quad (2.344)$$

[See Equation (2.324).] Hence, the mutual inductance can be written

$$M = \sqrt{L_1 L_2}. \quad (2.345)$$

Note that this result depends on the assumption that all of the magnetic flux produced by one solenoid passes through the other solenoid. In reality, some of the flux leaks out, so that the mutual inductance is somewhat less than that given in the previous formula. We can write

$$M = k \sqrt{L_1 L_2}, \quad (2.346)$$

where the dimensionless constant k is called the *coefficient of coupling*, and lies in the range $0 \leq k \leq 1$.

Suppose that the inner and outer solenoids have resistances R_1 and R_2 , respectively. If an instantaneous current I_1 flow through the inner solenoid then the voltage drop across it due to its resistance is $I_1 R_1$. The voltage drop due to the back-emf generated by the self inductance of the solenoid is $L_1 dI_1/dt$. (See Section 2.3.5.) There is also a back-emf due to inductive coupling with the outer solenoid. The magnetic flux through the inner solenoid due to the instantaneous current I_2 flowing through the outer solenoid is

$$\Phi_1 = M I_2. \quad (2.347)$$

Thus, by Faraday's law and Lenz's law, the back-emf induced in the inner solenoid is

$$\mathcal{E} = -\frac{d\Phi_1}{dt} = -M \frac{dI_2}{dt}. \quad (2.348)$$

[See Equation (2.284).] The voltage drop across the inner solenoid due to its mutual inductance with the top coil is minus this expression. Thus, the net voltage drop across the inner solenoid is

$$V_1 = R_1 I_1 + L_1 \frac{dI_1}{dt} + M \frac{dI_2}{dt}. \quad (2.349)$$

Likewise, the net voltage drop across the outer solenoid is

$$V_2 = R_2 I_2 + L_2 \frac{dI_2}{dt} + M \frac{dI_1}{dt}. \quad (2.350)$$

Suppose that, at time $t = 0$, we suddenly connect a battery of constant voltage V_1 to the inner solenoid. The outer solenoid is assumed to be open-circuited, or connected to a voltmeter of very high internal resistance, so that $I_2 = 0$. Because $I_2 = 0$, the circuit equation for the inner solenoid is

$$V_1 = R_1 I_1 + L_1 \frac{dI_1}{dt}, \quad (2.351)$$

where V_1 is constant, and $I_1(t = 0) = 0$. We have already seen the solution to this equation:

$$I_1 = \frac{V_1}{R_1} \left[1 - \exp\left(-\frac{R_1 t}{L_1}\right) \right]. \quad (2.352)$$

[See Equation (2.329).] The circuit equation for the outer solenoid is

$$V_2 = M \frac{dI_1}{dt}, \quad (2.353)$$

giving

$$V_2 = V_1 \frac{M}{L_1} \exp\left(-\frac{R_1 t}{L_1}\right). \quad (2.354)$$

It follows from Equation (2.346) that

$$V_2 = V_1 k \sqrt{\frac{L_2}{L_1}} \exp\left(-\frac{R_1 t}{L_1}\right). \quad (2.355)$$

Because $L_1/L_2 = N_1^2/N_2^2$ [see Equations (2.343) and (2.344)], we obtain

$$V_2 = V_1 k \frac{N_2}{N_1} \exp\left(-\frac{R_1 t}{L_1}\right). \quad (2.356)$$

Now,

$$\frac{V_2(t=0)}{V_1} = k \frac{N_2}{N_1}, \quad (2.357)$$

so if $N_2 \gg N_1$ then the voltage in the inner solenoid is considerably amplified in the outer solenoid. This effect is the basis for old-fashioned car ignition systems. A large voltage spike is induced in a secondary circuit (connected to a coil with very many turns) whenever the current in a primary circuit (connected to a coil with not so many turns) is either switched on or off. The primary circuit is connected to the car battery (whose voltage is typically 12 volts). The switching is done by a set of points, which are mechanically opened and closed as the engine turns. The large voltage spike induced in the secondary circuit, as the points are either opened or closed, causes a spark to jump across a gap in this circuit. This spark ignites a petrol/air mixture in one of the engine's cylinders. We might think that the optimum configuration is to have only one turn in the primary circuit, and many turns in the secondary circuit, so that the ratio N_2/N_1 is made as large as possible. However, this is not the case. Most of the magnetic flux generated by a single-turn primary coil is likely to miss the secondary coil altogether. This implies that the coefficient of coupling k is small, which reduces the voltage induced in the secondary circuit. Thus, we need a reasonable number of turns in the primary coil in order to localize the induced magnetic flux, so that it links effectively with the secondary coil.

2.3.8 Magnetic Energy

Suppose that, at $t = 0$, a solenoid of inductance L , and resistance R , is connected across the terminals of a battery of voltage V . The circuit equation is

$$V = L \frac{dI}{dt} + RI. \quad (2.358)$$

[See Equation (2.326).] The power output of the battery is VI . [Every charge q that goes around the circuit falls through a potential difference qV . In order to raise it back to the starting potential, so that it can perform another circuit, the battery must do work qV . See Section 2.1.5. The work done per unit time (i.e., the power) is nqV , where n is the number of charges per unit time passing a given point on the circuit. But, $I = nq$, so the power output is VI .] Thus, the net work done by the battery in raising the current in the circuit from zero at time $t = 0$ to I_T at time $t = T$ is

$$W = \int_0^T VI dt. \quad (2.359)$$

Using the circuit equation (2.358), we obtain

$$W = L \int_0^T I \frac{dI}{dt} dt + R \int_0^T I^2 dt, \quad (2.360)$$

giving

$$W = \frac{1}{2} L I_T^2 + R \int_0^T I^2 dt. \quad (2.361)$$

The second term on the right-hand side of the previous equation represents the irreversible conversion of electrical energy into heat energy by the resistor. (See Section 2.1.11.) The first term is the amount of energy stored in the solenoid at time T . This energy can be recovered after the solenoid is disconnected from the battery. Suppose that the battery is disconnected at time T . The circuit equation is now

$$0 = L \frac{dI}{dt} + RI, \quad (2.362)$$

giving

$$I = I_T \exp \left[-\frac{R}{L} (t - T) \right], \quad (2.363)$$

where we have made use of the initial condition $I(T) = I_T$. Thus, the current decays away exponentially. The energy stored in the solenoid is dissipated as heat in the resistor. The total heat energy appearing in the resistor after the battery is disconnected is

$$\int_T^\infty I^2 R dt = \frac{1}{2} L I_T^2, \quad (2.364)$$

where use has been made of Equation (2.363). Thus, the heat energy appearing in the resistor is equal to the energy stored in the solenoid. This energy is actually stored in the magnetic field generated inside the solenoid.

Consider, again, our circuit with two solenoids wound on top of one another. (See the previous section.) Suppose that each solenoid is connected to its own battery. The circuit equations are thus

$$V_1 = R_1 I_1 + L_1 \frac{dI_1}{dt} + M \frac{dI_2}{dt}, \quad (2.365)$$

$$V_2 = R_2 I_2 + L_2 \frac{dI_2}{dt} + M \frac{dI_1}{dt}, \quad (2.366)$$

where V_1 is the voltage of the battery in the first circuit, et cetera. The net work done by the two batteries in increasing the currents in the two circuits, from zero at time 0, to I_1 and I_2 at time T , respectively, is

$$\begin{aligned} W &= \int_0^T (V_1 I_1 + V_2 I_2) dt \\ &= \int_0^T (R_1 I_1^2 + R_2 I_2^2) dt + \frac{1}{2} L_1 I_1^2 + \frac{1}{2} L_2 I_2^2 \\ &\quad + M \int_0^T \left(I_1 \frac{dI_2}{dt} + I_2 \frac{dI_1}{dt} \right) dt. \end{aligned} \quad (2.367)$$

Thus,

$$W = \int_0^T (R_1 I_1^2 + R_2 I_2^2) dt + \frac{1}{2} L_1 I_1^2 + \frac{1}{2} L_2 I_2^2 + M I_1 I_2. \quad (2.368)$$

Clearly, the total magnetic energy stored in the two solenoids is

$$W_B = \frac{1}{2} L_1 I_1^2 + \frac{1}{2} L_2 I_2^2 + M I_1 I_2. \quad (2.369)$$

Note that the mutual inductance term increases the stored magnetic energy if I_1 and I_2 are of the same sign; that is, if the currents in the two solenoids flow in the same direction, so that they generate magnetic fields that reinforce one another. Conversely, the mutual inductance term decreases the stored magnetic energy if I_1 and I_2 are of the opposite sign. However, the total stored energy can never be negative, otherwise the coils would constitute a power source (a negative stored energy is equivalent to a positive generated energy). Thus,

$$\frac{1}{2} L_1 I_1^2 + \frac{1}{2} L_2 I_2^2 + M I_1 I_2 \geq 0, \quad (2.370)$$

which can be written

$$\frac{1}{2} (\sqrt{L_1} I_1 + \sqrt{L_2} I_2)^2 - I_1 I_2 (\sqrt{L_1 L_2} - M) \geq 0, \quad (2.371)$$

assuming that $I_1 I_2 < 0$. It follows that

$$M \leq \sqrt{L_1 L_2}. \quad (2.372)$$

The equality sign corresponds to the situation in which all of the magnetic flux generated by one solenoid passes through the other. If some of the flux misses then the inequality sign is appropriate. In fact, the previous formula is valid for any two inductively coupled circuits, and effectively sets an upper limit on their mutual inductance.

We intimated previously that the energy stored in an solenoid is actually stored in the surrounding magnetic field. Let us now obtain an explicit formula for the energy stored in a magnetic field. Consider an ideal cylindrical solenoid. The energy stored in the solenoid when a current I flows through it is

$$W = \frac{1}{2} L I^2, \quad (2.373)$$

where L is the self inductance. We know that

$$L = \mu_0 N^2 \pi r^2 l, \quad (2.374)$$

where N is the number of turns per unit length of the solenoid, r the radius, and l the length. [See Equation (2.324).] The magnetic field inside the solenoid is approximately uniform, with magnitude

$$B = \mu_0 N I, \quad (2.375)$$

and is approximately zero outside the solenoid. [See Equation (2.279).] Equation (2.373) can be rewritten

$$W = \frac{B^2}{2\mu_0} \mathcal{V}, \quad (2.376)$$

where $\mathcal{V} = \pi r^2 l$ is the volume of the solenoid. The previous formula strongly suggests that a magnetic field possesses an energy density

$$U = \frac{B^2}{2\mu_0}. \quad (2.377)$$

Let us now examine a more general proof of the previous formula. Consider a system of N circuits (labeled $i = 1$ to N), each carrying a current I_i . The magnetic flux through the i th circuit is written [see Equation (2.315)]

$$\Phi_i = \int \mathbf{B} \cdot d\mathbf{S}_i = \oint \mathbf{A} \cdot d\mathbf{r}_i, \quad (2.378)$$

where $\mathbf{B} = \nabla \times \mathbf{A}$, and $d\mathbf{S}_i$ and $d\mathbf{r}_i$ denote a surface element and a line element of this circuit, respectively. The back-emf induced in the i th circuit follows from Faraday's law:

$$\mathcal{E}_i = -\frac{d\Phi_i}{dt}. \quad (2.379)$$

[See Equation (2.284).] The rate of work of the battery that maintains the current I_i in the i th circuit against this back-emf is

$$P_i = -I_i \mathcal{E}_i = I_i \frac{d\Phi_i}{dt}. \quad (2.380)$$

Thus, the total work required to raise the currents in the N circuits from zero at time 0, to I_{0i} at time T , is

$$W = \sum_{i=1,N} \int_0^T I_i \frac{d\Phi_i}{dt} dt. \quad (2.381)$$

The previous expression for the work done is, of course, equivalent to the total energy stored in the magnetic field surrounding the various circuits. This energy is independent of the manner in which the currents are set up. Suppose, for the sake of simplicity, that the currents are ramped up linearly, so that

$$I_i = I_{0i} \frac{t}{T}. \quad (2.382)$$

The fluxes are proportional to the currents, so they must also ramp up linearly; that is,

$$\Phi_i = \Phi_{0i} \frac{t}{T}. \quad (2.383)$$

It follows that

$$W = \sum_{i=1,N} \int_0^T I_{0i} \Phi_{0i} \frac{t}{T^2} dt, \quad (2.384)$$

giving

$$W = \frac{1}{2} \sum_{i=1, N} I_{0i} \Phi_{0i}. \quad (2.385)$$

So, if instantaneous currents I_i flow in the N circuits, which link instantaneous fluxes Φ_i , then the instantaneous stored energy is

$$W = \frac{1}{2} \sum_{i=1, N} I_i \Phi_i. \quad (2.386)$$

Equations (2.378) and (2.386) imply that

$$W = \frac{1}{2} \sum_{i=1, N} I_i \oint \mathbf{A} \cdot d\mathbf{r}_i. \quad (2.387)$$

It is convenient, at this stage, to replace our N line currents by N current distributions of small, but finite, cross-sectional area. Equation (2.387) transforms to give

$$W = \frac{1}{2} \int_V \mathbf{A} \cdot \mathbf{j} dV, \quad (2.388)$$

where V is a volume that contains all of the circuits. Note that for an element of the i th circuit, $\mathbf{j} = I_i d\mathbf{r}_i / (dr_i A_i)$ and $dV = dr_i A_i$, where A_i is the cross-sectional area of the circuit. Now, $\mu_0 \mathbf{j} = \nabla \times \mathbf{B}$ [see Equation (2.271)], so

$$W = \frac{1}{2\mu_0} \int_V \mathbf{A} \cdot \nabla \times \mathbf{B} dV. \quad (2.389)$$

However,

$$\nabla \cdot (\mathbf{A} \times \mathbf{B}) \equiv \mathbf{B} \cdot \nabla \times \mathbf{A} - \mathbf{A} \cdot \nabla \times \mathbf{B} \quad (2.390)$$

(see Section A.24), which implies that

$$W = \frac{1}{2\mu_0} \int_V [-\nabla \cdot (\mathbf{A} \times \mathbf{B}) + \mathbf{B} \cdot \nabla \times \mathbf{A}] dV. \quad (2.391)$$

Using the divergence theorem (see Section A.20), and $\mathbf{B} = \nabla \times \mathbf{A}$, we obtain

$$W = -\frac{1}{2\mu_0} \oint_S \mathbf{A} \times \mathbf{B} \cdot d\mathbf{S} + \frac{1}{2\mu_0} \int_V B^2 dV, \quad (2.392)$$

where S is the bounding surface of some volume V . Let us take this surface to infinity. It is easily demonstrated that the magnetic field generated by a current loop falls off like r^{-3} at large distances. (See Section 2.2.7.) The vector potential falls off like r^{-2} . However, the area of surface S only increases like r^2 . It follows that the surface integral is negligible in the limit $r \rightarrow \infty$. Thus, the previous expression reduces to

$$W = \int \frac{B^2}{2\mu_0} dV, \quad (2.393)$$

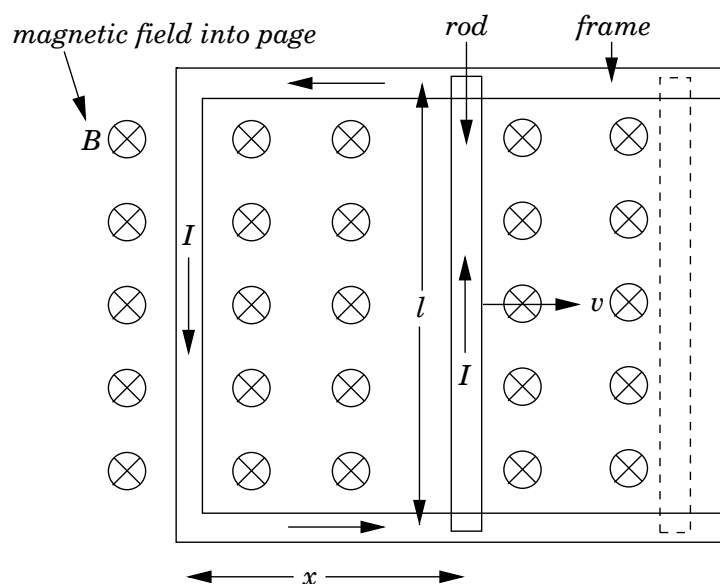


Figure 2.31: Motional emf.

where the integral is over all space. Because this expression is valid for any magnetic field whatsoever, we can safely conclude that the energy density of a general magnetic field generated by a system of electrical circuits is given by

$$U = \frac{B^2}{2\mu_0}. \quad (2.394)$$

Note, finally, that the fact that a magnetic field possesses an energy density demonstrates that it has a real physical existence, and is not merely an aid to calculating the forces that current-carrying wires exert on one another.

2.3.9 Motional Emf

Consider a simple circuit in which a conducting rod of length l slides along a U-shaped conducting frame in the presence of a uniform magnetic field. This circuit is illustrated in Figure 2.31. Suppose, for the sake of simplicity, that the magnetic field is directed perpendicular to the plane of the circuit. To be more exact, the magnetic field is directed into the page in the figure. Suppose, further, that the rod moves to the right at the constant speed v .

The magnetic flux passing through the circuit is simply the product of the perpendicular magnetic field-strength, B , and the area of the circuit, lx , where x determines the position of the sliding rod. Thus,

$$\Phi_B = Blx. \quad (2.395)$$

Now, the rod moves a distance $dx = v dt$ in a time interval dt , so in the same time interval the magnetic flux passing through the circuit increases by

$$d\Phi_B = B l dx = B l v dt. \quad (2.396)$$

It follows, from Faraday's law [see Equation (2.284)], that the magnitude of the emf V generated around the circuit is given by

$$V = \frac{d\Phi_B}{dt} = Blv. \quad (2.397)$$

Thus, the emf generated in the circuit by the moving rod is simply the product of the magnetic field-strength, the length of the rod, and the velocity of the rod. If the magnetic field is not perpendicular to the circuit, but instead subtends an angle θ with respect to the normal direction to the plane of the circuit, then it is easily demonstrated that the so-called *motional emf* generated in the circuit by the moving rod is

$$V = B_{\perp} lv, \quad (2.398)$$

where $B_{\perp} = B \cos \theta$ is the component of the magnetic field that is perpendicular to the plane of the circuit.

Because the magnetic flux linking the circuit increases in time, by Lenz's law, the emf acts in the negative direction (i.e., in the opposite sense to the fingers of a right-hand, if the thumb points along the direction of the magnetic field). The emf, V , therefore, acts in a counter-clockwise direction in the figure. If R is the total resistance of the circuit then this emf drives an counter-clockwise electric current of magnitude $I = V/R$ around the circuit. Of course, this current generates a magnetic field that acts to reduce the increase in the magnetic flux passing through the circuit.

But, where does the motional emf come from? Let us again remind ourselves what an we mean by an emf. When we say that an emf V acts around the circuit in the counter-clockwise direction, what we really mean is that a charge q that circulates once around the circuit in a counter-clockwise direction acquires the energy qV . The only manner in which the charge can acquire this energy is if something does work on it as it circulates. Let us assume that the charge circulates very slowly. The magnetic field exerts a negligibly small force on the charge when it is traversing the non-moving part of the circuit (because the charge is moving very slowly). However, when the charge is traversing the moving rod it experiences an upward (in the figure) magnetic force of magnitude $f = qvB$ (assuming that $q > 0$). (See Section 2.2.4.) The net work done on the charge by this force as it traverses the rod is

$$W' = qvBl = qV, \quad (2.399)$$

because $V = Blv$. Thus, it would appear that the motional emf generated around the circuit can be accounted for in terms of the magnetic force exerted on charges traversing the moving rod.

However, there is something seriously wrong with the previous explanation. We seem to be saying that the charge acquires the energy qV from the magnetic field as it moves around the circuit once in a counter-clockwise direction. But, this is impossible, because a magnetic field cannot do work on an electric charge. (See Section 2.2.4.)

Let us look at the problem from the point of view of a charge q traversing the moving rod. In the frame of reference of the rod, the charge only moves very slowly, so the magnetic force acting on it is negligible. In fact, only an electric field can exert a significant force on a slowly moving charge. In order to account for the motional emf generated around the circuit, we need the charge to experience an upward force of magnitude qvB . The only way in which this is possible is if the charge sees an upward pointing electric field of magnitude

$$E = vB. \quad (2.400)$$

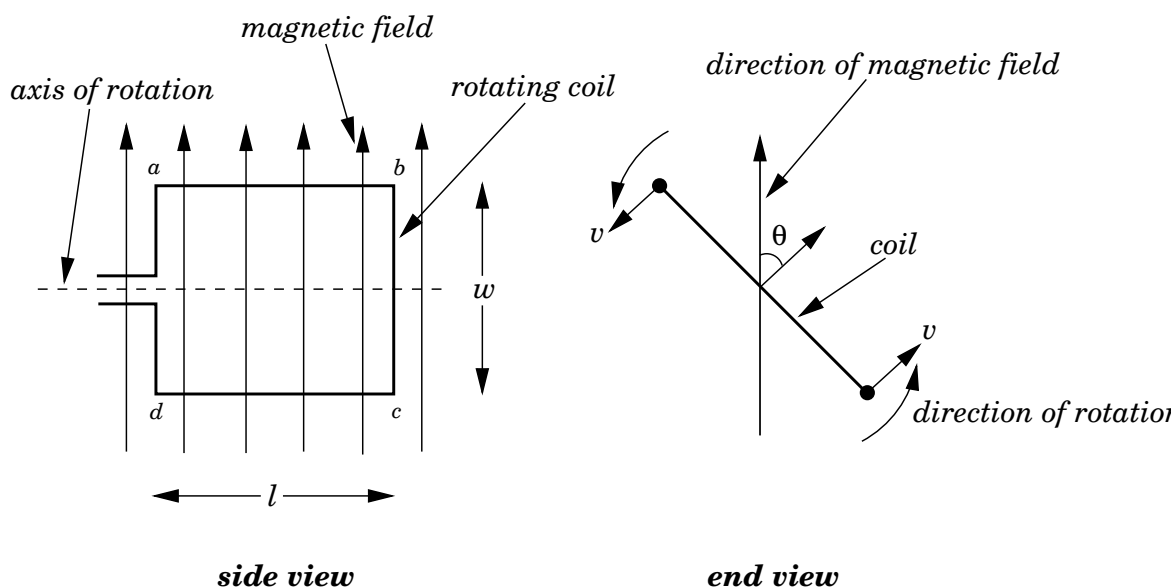


Figure 2.32: An alternating current generator.

In other words, although there is no electric field in the laboratory frame, there is an electric field in the frame of reference of the moving rod, and it is this field that does the necessary amount of work on charges moving around the circuit in order to account for the existence of the motional emf, $V = Blv$.

More generally, if a conductor moves in the laboratory frame with velocity \mathbf{v} in the presence of a magnetic field \mathbf{B} then a charge q inside the conductor experiences a magnetic force $\mathbf{f} = q \mathbf{v} \times \mathbf{B}$. In the frame of the conductor, in which the charge is essentially stationary, the same force takes the form of an electric force $\mathbf{f} = q \mathbf{E}$, where \mathbf{E} is the electric field in the frame of reference of the conductor. Thus, if a conductor moves with velocity \mathbf{v} through a magnetic field \mathbf{B} then the electric field \mathbf{E} that appears in the rest frame of the conductor is given by

$$\mathbf{E} = \mathbf{v} \times \mathbf{B}. \quad (2.401)$$

(See Section 3.4.1.) This electric field is the ultimate origin of the motional emfs that are generated whenever circuits move with respect to magnetic fields.

2.3.10 Alternating Current Generators

An *electric generator*, or *dynamo*, is a device that converts mechanical energy into electrical energy. The simplest practical generator consists of a rectangular coil rotating in a uniform magnetic field. The magnetic field is usually supplied by a permanent magnet. This setup is illustrated in Figure 2.32.

Let l be the length of the coil along its axis of rotation, and w the width of the coil perpendicular to this axis. Suppose that the coil rotates at constant angular velocity ω in a uniform magnetic field of strength B . The velocity v with which the two long sides of the coil (i.e., sides ab and cd)

move through the magnetic field is simply the product of the angular velocity of rotation ω and the distance $w/2$ of each side from the axis of rotation, so $v = \omega w/2$. The motional emf induced in each side is given by $V = B_{\perp} l v$, where B_{\perp} is the component of the magnetic field perpendicular to instantaneous direction of motion of the side in question. If the direction of the magnetic field subtends an angle θ with the normal direction to the coil, as shown in the figure, then $B_{\perp} = B \sin \theta$. Thus, the magnitude of the motional emf generated in sides ab and cd is

$$V_{ab} = \frac{B w l \omega \sin \theta}{2} = \frac{B A \omega \sin \theta}{2}, \quad (2.402)$$

where $A = w l$ is the area of the coil. The emf is zero when $\theta = 0^{\circ}$ or 180° , because the direction of motion of sides ab and cd is parallel to the direction of the magnetic field in these cases. The emf attains its maximum value when $\theta = 90^{\circ}$ or 270° , because the direction of motion of sides ab and cd is perpendicular to the direction of the magnetic field in these cases. Incidentally, it is clear, from symmetry, that no net motional emf is generated in sides bc and da of the coil.

Suppose that the direction of rotation of the coil is such that side ab is moving into the page in Figure 2.32 (side view), whereas side cd is moving out of the page. The motional emf induced in side ab acts from a to b . Likewise, the motional emf induce in side cd acts from c to d . It can be seen that both emfs act in the clockwise direction around the coil. [The direction of the emf is the same as the direction of the electric field seen in the rest frame of the sides. See Equation (2.401).] Thus, the net emf V acting around the coil is $2 V_{ab}$. If the coil has N turns then the net emf becomes $2 N V_{ab}$. Hence, the general expression for the emf generated around a steadily-rotating, multi-turn coil in a uniform magnetic field is

$$V = N B A \omega \sin(\omega t), \quad (2.403)$$

where we have written $\theta = \omega t$ for a steadily rotating coil (assuming that $\theta = 0$ at $t = 0$). This expression can also be written

$$V = V_{\max} \sin(2\pi f t), \quad (2.404)$$

where

$$V_{\max} = 2\pi N B A f \quad (2.405)$$

is the peak emf produced by the generator, and $f = \omega/2\pi$ is the number of complete rotations the coils executes per second. Thus, the peak emf is directly proportional to the area of the coil, the number of turns in the coil, the rotation frequency of the coil, and the magnetic field-strength.

Figure 2.33 shows the emf specified in Equation (2.404) plotted as a function of time. It can be seen that the variation of the emf with time is sinusoidal in nature. The emf attains its peak values when the plane of the coil is parallel to the plane of the magnetic field, passes through zero when the plane of the coil is perpendicular to the magnetic field, and reverses sign every half period of revolution of the coil. The emf is periodic (i.e., it continually repeats the same pattern in time), with period $T = 1/f$ (which is, of course, the rotation period of the coil).

Suppose that some electrical load (e.g., a light-bulb, or an electric heating element) of resistance R is connected across the terminals of the generator. In practice, this is achieved by connecting the

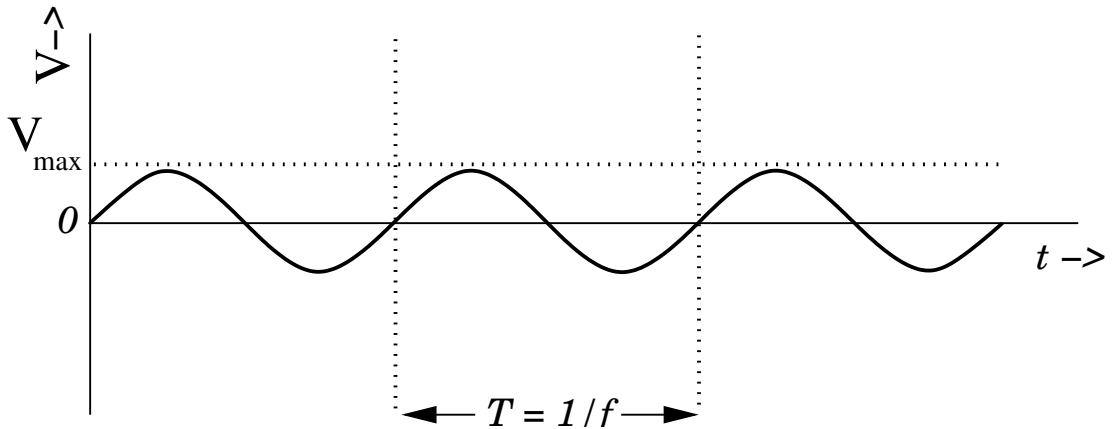


Figure 2.33: Emf generated by a steadily rotating AC generator.

two ends of the coil to rotating rings that are then connected to the external circuit by means of metal brushes. According to Ohm's law, the current I that flows in the load is given by

$$I = \frac{V}{R} = \frac{V_{\max}}{R} \sin(2\pi f t). \quad (2.406)$$

(See Section 2.1.11.) Note that this current is constantly changing direction, just like the emf of the generator. Hence, the type of generator described previously is usually termed an *alternating current*, or AC, generator.

The current I that flows through the load must also flow around the coil. Because the coil is situated in a magnetic field, this current gives rise to a torque acting on the coil which, as is easily demonstrated, acts to slow down its rotation. Suppose, as before, that side ab is moving into the page in Figure 2.32 (side view), whereas side cd is moving out of the page, and the current is circulating in a clockwise sense. Side ab experiences a magnetic force per unit length $F_{ab} = \mathbf{I} \times \mathbf{B} = IB \sin \theta$ that acts to oppose its motion. (See Section 2.2.2.) Hence, the braking force acting on the side is $f_{ab} = F_{ab} l = IB \sin \theta l$. Thus, the braking torque acting on the side is $\tau_{ab} = f_{ab} w/2 = IB \sin \theta l w/2 = IB \sin \theta A/2$, where $A = lw$ is the area of the coil. Side cd experiences an equal torque. So, taking into account the fact that the coils has N turns, the net braking torque τ acting on the coil is given by

$$\tau = N I B A \sin \theta. \quad (2.407)$$

It follows from Equation (2.403) that

$$\tau = \frac{VI}{\omega}, \quad (2.408)$$

because $V = N B A \omega \sin \theta$. An external torque that is equal and opposite to the braking torque must be applied to the coil if it is to rotate uniformly, as was initially assumed above. The rate P at which this external torque does work is equal to the product of the torque τ and the angular velocity ω of the coil. (See Section 1.7.4.) Thus,

$$P = \tau \omega = VI. \quad (2.409)$$

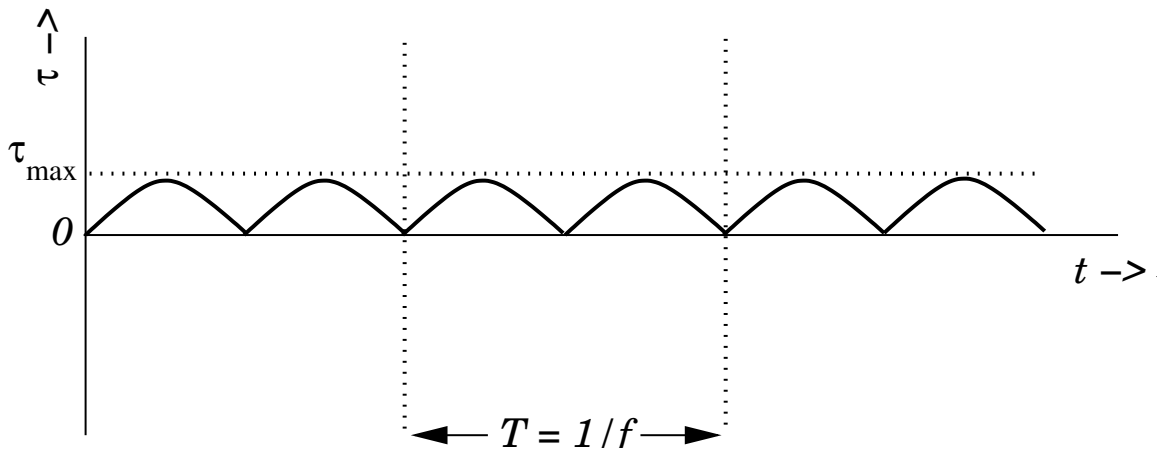


Figure 2.34: The braking torque in a steadily rotating AC generator.

Not surprisingly, the rate at which the external torque performs work exactly matches the rate $V I$ at which electrical energy is generated in the circuit comprising the rotating coil and the load.

Equations (2.403), (2.406), and (2.408) yield

$$\tau = \tau_{\max} \sin^2(2\pi f t), \quad (2.410)$$

where $\tau_{\max} = (V_{\max})^2 / (2\pi f R)$. Figure 2.34 shows the braking torque τ plotted as a function of time t , according to Equation (2.410). It can be seen that the torque is always of the same sign (i.e., it always acts in the same direction, so as to continually oppose the rotation of the coil), but is not constant in time. Instead, it pulsates periodically with period T . The braking torque attains its maximum value whenever the plane of the coil is parallel to the plane of the magnetic field, and is zero whenever the plane of the coil is perpendicular to the magnetic field. It is clear that the external torque needed to keep the coil rotating at a constant angular velocity must also pulsate in time with period T . A constant external torque would give rise to a non-uniformly rotating coil, and, hence, to an alternating emf that varies with time in a more complicated manner than $\sin(2\pi f t)$.

Virtually all commercial power stations generate electricity using AC generators. The external power needed to turn the generating coil is usually supplied by a steam turbine (steam blasting against fan-like blades that are forced into rotation). Water is vaporized to produce high pressure steam by burning coal, or by using the energy released inside a nuclear reactor. Of course, in hydroelectric power stations, the power needed to turn the generator coil is supplied by a water turbine (which is similar to a steam turbine, except that falling water plays the role of the steam). More recently, a new type of power station has been developed in which the power needed to rotate the generating coil is supplied by a gas turbine (basically, a large jet engine that burns natural gas). In the U.S. and Canada, the alternating electrical signal generated by power stations and fed into ordinary households, which is known as *mains electricity*, oscillates at $f = 60\text{Hz}$, which implies that the generator coils in power stations rotate exactly sixty times a second. In Europe and Asia, the oscillation frequency of mains electricity is $f = 50\text{Hz}$.

2.3.11 Alternating Current Circuits

Alternating current (AC) circuits are made up of voltage sources and three different types of passive elements. These are resistors, inductors (i.e., small solenoids), and capacitors. Resistors satisfy Ohm's law,

$$V = IR, \quad (2.411)$$

where R is the resistance, I the current flowing through the resistor, and V the voltage drop across the resistor (in the direction in which the current flows). (See Section 2.1.11.) Inductors satisfy

$$V = L \frac{dI}{dt}, \quad (2.412)$$

where L is the inductance. [See Equation (2.325).] Finally, capacitors obey

$$V = \frac{q}{C} = \int_0^t I dt / C, \quad (2.413)$$

where C is the capacitance, q is the charge stored on the plate with the most positive potential, and $I = 0$ for $t < 0$. (See Section 2.1.13.) Note that any passive component of a real electrical circuit can always be represented as a combination of ideal resistors, inductors, and capacitors.

Let us consider the classic LCR circuit, which consists of an inductor, L , a capacitor, C , and a resistor, R , all connected in series with an voltage source, V . See Figure 2.35. The circuit equation is obtained by setting the input voltage V equal to the sum of the voltage drops across the three passive elements in the circuit. Thus,

$$V = IR + L \frac{dI}{dt} + \int_0^t I dt / C. \quad (2.414)$$

This is an integro-differential equation which, in general, is quite difficult to solve. Suppose, however, that both the voltage and the current oscillate at some fixed angular frequency, ω , so that

$$V(t) = V_0 \exp(i \omega t), \quad (2.415)$$

$$I(t) = I_0 \exp(i \omega t), \quad (2.416)$$

where $i = \sqrt{-1}$, and the physical solution is understood to be the real part of the previous expressions. The assumed behavior of the voltage and current is clearly relevant to electrical circuits powered by mains electricity (which oscillates at 60 hertz in the U.S. and Canada).

Equations (2.414)–(2.416) yield

$$V_0 \exp(i \omega t) = I_0 \exp(i \omega t) R + L i \omega I_0 \exp(i \omega t) + \frac{I_0 \exp(i \omega t)}{i \omega C}, \quad (2.417)$$

giving

$$V_0 = I_0 \left(i \omega L + \frac{1}{i \omega C} + R \right). \quad (2.418)$$

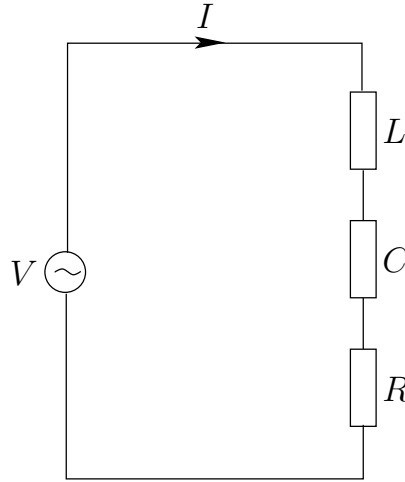


Figure 2.35: An LCR circuit.

It is helpful to define the *impedance* of the circuit:

$$Z = \frac{V}{I} = i\omega L + \frac{1}{i\omega C} + R. \quad (2.419)$$

Impedance is a generalization of the concept of resistance. In general, the impedance of an AC circuit is a complex quantity.

The average power output of the voltage source is

$$P = \langle V(t) I(t) \rangle, \quad (2.420)$$

where the average is taken over one period of the oscillation. Let us, first of all, calculate the power using real (rather than complex) voltages and currents. We can write

$$V(t) = |V_0| \cos(\omega t), \quad (2.421)$$

$$I(t) = |I_0| \cos(\omega t - \theta), \quad (2.422)$$

where θ is the phase-lag of the current with respect to the voltage. It follows that

$$\begin{aligned} P &= |V_0| |I_0| \int_{\omega t=0}^{\omega t=2\pi} \cos(\omega t) \cos(\omega t - \theta) \frac{d(\omega t)}{2\pi} \\ &= |V_0| |I_0| \int_{\omega t=0}^{\omega t=2\pi} \cos(\omega t) [\cos(\omega t) \cos \theta + \sin(\omega t) \sin \theta] \frac{d(\omega t)}{2\pi}, \end{aligned} \quad (2.423)$$

giving

$$P = \frac{1}{2} |V_0| |I_0| \cos \theta, \quad (2.424)$$

because $\langle \cos(\omega t) \sin(\omega t) \rangle = 0$ and $\langle \cos(\omega t) \cos(\omega t) \rangle = 1/2$. Here, $\langle \cdots \rangle \equiv \int_{\omega t=0}^{\omega t=2\pi} (\cdots) d(\omega t)/(2\pi)$. In complex representation, the voltage and the current are written

$$V(t) = |V_0| \exp(i\omega t), \quad (2.425)$$

$$I(t) = |I_0| \exp[i(\omega t - \theta)]. \quad (2.426)$$

Now,

$$\frac{1}{2} (VI^* + V^*I) = |V_0||I_0| \cos \theta. \quad (2.427)$$

It follows from Equation (2.424) that

$$P = \frac{1}{4} (VI^* + V^*I) = \frac{1}{2} \operatorname{Re}(VI^*). \quad (2.428)$$

Making use of Equation (2.419), we find that

$$P = \frac{1}{2} \operatorname{Re}(Z) |I|^2 = \frac{1}{2} \frac{\operatorname{Re}(Z) |V|^2}{|Z|^2}. \quad (2.429)$$

Note that power dissipation is associated with the real part of the impedance. For the specific case of an LCR circuit,

$$P = \frac{1}{2} R |I_0|^2. \quad (2.430)$$

[See Equation (2.419).] We conclude that only the resistor dissipates energy in this circuit. The inductor and the capacitor both store energy, but they eventually return it to the circuit without dissipation.

According to Equation (2.419), the amplitude of the current that flows in an LCR circuit, for a given amplitude of the input voltage, is given by

$$|I_0| = \frac{|V_0|}{|Z|} = \frac{|V_0|}{\sqrt{(\omega L - 1/\omega C)^2 + R^2}}. \quad (2.431)$$

As can be seen from Figure 2.36, the response of the circuit is resonant, peaking at $\omega = 1/\sqrt{LC}$, and reaching $1/\sqrt{2}$ of the peak value at $\omega = 1/\sqrt{LC} \pm R/(2L)$ (assuming that $R \ll \sqrt{L/C}$). For this reason, LCR circuits are used in analog radio tuners to filter out signals whose frequencies fall outside a given band.

The phase-lag of the current with respect to the voltage is given by

$$\theta = \arg(Z) = \tan^{-1} \left(\frac{\omega L - 1/\omega C}{R} \right). \quad (2.432)$$

[See Equation (2.419).] As can be seen from Figure 2.36, the phase-lag varies from -90° for frequencies significantly below the resonant frequency, to zero at the resonant frequency ($\omega = 1/\sqrt{LC}$), to $+90^\circ$ for frequencies significantly above the resonant frequency.

It is clear that, in conventional AC circuits, the circuit equation reduces to a simple algebraic equation, and that the behavior of the circuit is summed up by the complex impedance, Z . The real part of Z tells us the power dissipated in the circuit, the magnitude of Z gives the ratio of the peak current to the peak voltage, and the argument of Z gives the phase-lag of the current with respect to the voltage.

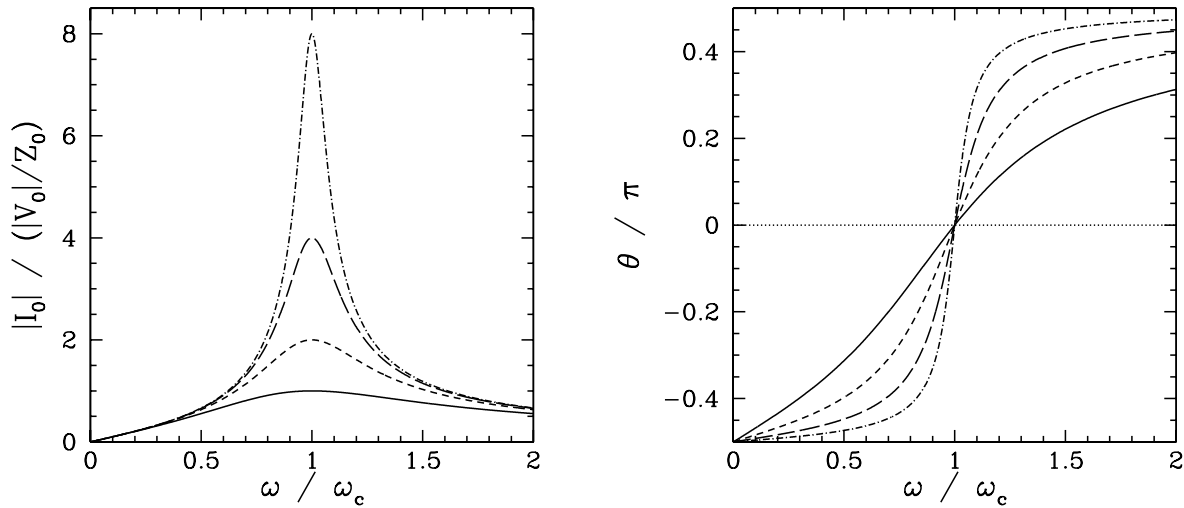


Figure 2.36: The characteristics of an LCR circuit. The left-hand and right-hand panes show the amplitude and phase-lag of the current versus frequency, respectively. Here, $\omega_c = 1/\sqrt{LC}$ and $Z_0 = \sqrt{L/C}$. The solid, short-dashed, long-dashed, and dot-dashed curves correspond to $R/Z_0 = 1, 1/2, 1/4,$ and $1/8,$ respectively.

2.3.12 Alternating Current Motors

The first electric dynamo was constructed in 1831 by Michael Faraday. An electric dynamo is, of course, a device that transforms mechanical energy into electrical energy. An *electric motor*, on the other hand, is a device that transforms electrical energy into mechanical energy. In other words, an electric motor is an electric dynamo run in reverse. It took a surprisingly long time for scientists in the nineteenth century to realize this. In fact, the message only really sank home after a fortuitous accident during the 1873 Vienna World Exposition. A large hall was filled with modern gadgets. One of these gadgets, a steam engine driven dynamo, was producing electric power when a workman unwittingly connected the output leads from another dynamo to the energized circuit. Almost immediately, the latter dynamo started to whirl around at great speed. The dynamo was, in effect, transformed into an electric motor.

An AC electric motor consists of the same basic elements as an AC electric generator; that is, a multi-turn coil that is free to rotate in a constant magnetic field. Furthermore, the rotating coil is connected to the external circuit in just the same manner as in an AC generator; that is, via two slip-rings attached to metal brushes. Suppose that an external voltage source of emf V is connected across the motor. It is assumed that V is an alternating emf, so that

$$V = V_{\max} \sin(2\pi f t), \quad (2.433)$$

where V_{\max} is the peak voltage, and f is the alternation frequency. Such an emf could be obtained from mains electricity. In this case, $V_{\max} = 110\text{V}$ and $f = 60\text{Hz}$ in the U.S. and Canada, whereas $V_{\max} = 220\text{V}$ and $f = 50\text{Hz}$ in Europe and Asia. The external emf drives an alternating current

$$I = I_{\max} \sin(2\pi f t) \quad (2.434)$$

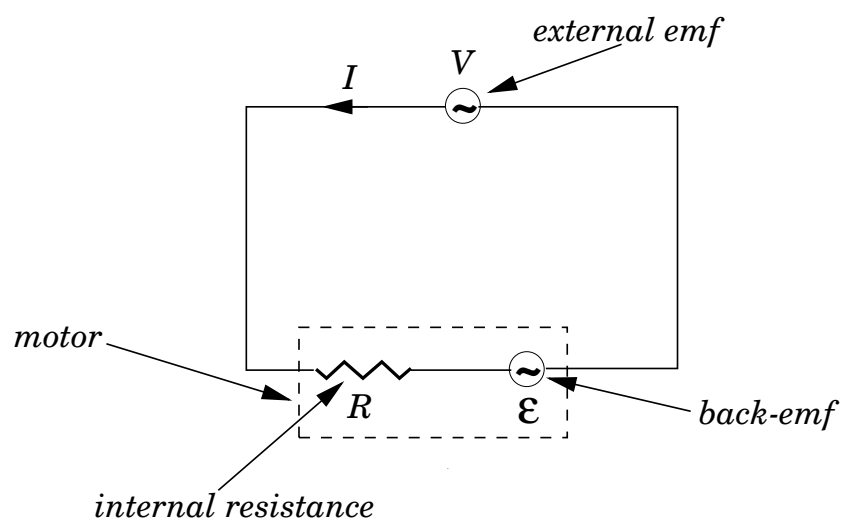


Figure 2.37: Circuit diagram for an AC motor connected to an external AC emf source.

around the external circuit, and through the motor. As this current flows around the coil, the magnetic field exerts a torque on the coil, which causes it to rotate. The motor eventually attains a steady state in which the rotation frequency of the coil matches the alternation frequency of the external emf. In other words, the steady-state rotation frequency of the coil is f . Now a coil rotating in a magnetic field generates an emf, \mathcal{E} . It is easily demonstrated that this emf acts to oppose the circulation of the current around the coil; that is, the induced emf acts in the opposite direction to the external emf. For an N -turn coil of cross-sectional area A , rotating with frequency f in a magnetic field B , the back-emf, \mathcal{E} , is given by

$$\mathcal{E} = \mathcal{E}_{\max} \sin(2\pi f t), \quad (2.435)$$

where

$$\mathcal{E}_{\max} = 2\pi N B A f, \quad (2.436)$$

and use has been made of the results of Section 2.3.10.

Figure 2.37 shows the circuit in question. A circle with a wavy line inside is the conventional way of indicating an AC voltage source. The motor is modeled as a resistor R , that represents the internal resistance of the motor, in series with the back-emf, \mathcal{E} . Of course, the back-emf acts in the opposite direction to the external emf, V . Application of Ohm's law (see Section 2.1.11) around the circuit gives

$$V = IR + \mathcal{E}, \quad (2.437)$$

or

$$V_{\max} \sin(2\pi f t) = I_{\max} R \sin(2\pi f t) + \mathcal{E}_{\max} \sin(2\pi f t), \quad (2.438)$$

which reduces to

$$V_{\max} = I_{\max} R + \mathcal{E}_{\max}. \quad (2.439)$$

The rate P at which the motor gains electrical energy from the external circuit is given by

$$P = \mathcal{E}I = P_{\max} \sin^2(2\pi f t), \quad (2.440)$$

where

$$P_{\max} = \mathcal{E}_{\max} I_{\max} = \frac{\mathcal{E}_{\max} (V_{\max} - \mathcal{E}_{\max})}{R}. \quad (2.441)$$

By conservation of energy, P is also the rate at which the motor performs mechanical work. Note that the rate at which the motor does mechanical work is not constant in time, but, instead, pulsates at the rotation frequency of the coil. It is possible to construct a motor that performs work at a more uniform rate by employing more than one coil rotating about the same axis.

As long as $V_{\max} > \mathcal{E}_{\max}$, the rate at which the motor performs mechanical work is positive (i.e., the motor does useful work). However, if $V_{\max} < \mathcal{E}_{\max}$ then the rate at which the motor performs work becomes negative. This means that we must perform mechanical work on the motor in order to keep it rotating, which is another way of saying that the motor does not perform useful work. Clearly, in order for an AC motor to perform useful work, the external emf, V , must be able to overcome the back-emf, \mathcal{E} , induced in the motor (i.e., $V_{\max} > \mathcal{E}_{\max}$).

2.3.13 Transformers

A *transformer* is a device for stepping-up, or stepping-down, the voltage of an alternating electric signal. Without efficient transformers, the transmission and distribution of AC electric power over long distances would be impossible. Figure 2.38 shows the circuit diagram of a typical transformer. There are two circuits. Namely, the primary circuit, and the secondary circuit. There is no direct electrical connection between the two circuits, but each circuit contains a coil that links it inductively to the other circuit. In real transformers, the two coils are wound onto the same iron core. The purpose of the iron core is to channel the magnetic flux generated by the current flowing around the primary coil, so that as much of it as possible also links the secondary coil. The common magnetic flux linking the two coils is conventionally denoted in circuit diagrams by a number of parallel straight-lines drawn between the coils.

Let us consider a particularly simple transformer in which the primary and secondary coils are solenoids that share the same air-filled core. Suppose that l is the length of the core, and A is its cross-sectional area. Let N_1 be the total number of turns in the primary coil, and let N_2 be the total number of turns in the secondary coil. Suppose that an alternating voltage

$$v_1 = V_1 \cos(\omega t) \quad (2.442)$$

is fed into the primary circuit from some external AC power source. Here, V_1 is the peak voltage in the primary circuit, and ω is the alternation frequency (in radians per second). The current driven around the primary circuit is written

$$i_1 = I_1 \sin(\omega t), \quad (2.443)$$

where I_1 is the peak current. This current generates a changing magnetic flux, in the core of the solenoid, that links the secondary coil, and, thereby, inductively generates the alternating emf

$$v_2 = V_2 \cos(\omega t) \quad (2.444)$$

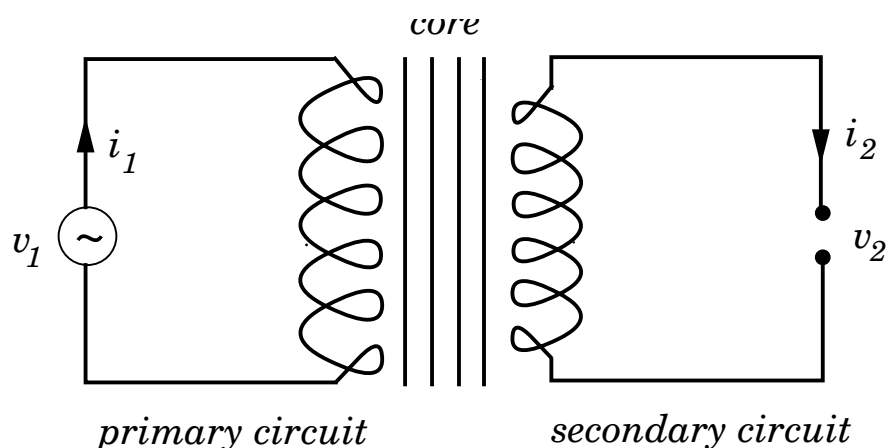


Figure 2.38: Circuit diagram of a transformer.

in the secondary circuit, where V_2 is the peak voltage. Suppose that this emf drives an alternating current

$$i_2 = I_2 \sin(\omega t) \quad (2.445)$$

around the secondary circuit, where I_2 is the peak current.

The circuit equation for the primary circuit is written

$$v_1 - L_1 \frac{di_1}{dt} - M \frac{di_2}{dt} = 0, \quad (2.446)$$

assuming that there is negligible resistance in this circuit. The first term in the previous equation is the externally generated emf. The second term is the back-emf due to the self inductance L_1 of the primary coil. (See Section 2.3.5.) The final term is the emf due to the mutual inductance M of the primary and secondary coils. (See Section 2.3.7.) In the absence of any significant resistance in the primary circuit, these three emfs must add up to zero. Equations (2.442)–(2.446) can be combined to give

$$V_1 = \omega (L_1 I_1 + M I_2), \quad (2.447)$$

because

$$\frac{d \sin(\omega t)}{dt} = \omega \cos(\omega t). \quad (2.448)$$

The alternating emf generated in the secondary circuit consists of the emf generated by the self inductance, L_2 , of the secondary coil, plus the emf generated by the mutual inductance of the primary and secondary coils. Thus,

$$v_2 = L_2 \frac{di_2}{dt} + M \frac{di_1}{dt}. \quad (2.449)$$

Equations (2.443)–(2.445), (2.448), and (2.449) yield

$$V_2 = \omega (L_2 I_2 + M I_1). \quad (2.450)$$

Now, the instantaneous power output of the external AC power source that drives the primary circuit is

$$P_1 = i_1 v_1. \quad (2.451)$$

Likewise, the instantaneous electrical energy per unit time transferred inductively from the primary to the secondary circuit is

$$P_2 = i_2 v_2. \quad (2.452)$$

If resistive losses in the primary and secondary circuits are negligible, as is assumed to be the case, then, by energy conservation, these two powers must equal one another at all times. Thus,

$$i_1 v_1 = i_2 v_2, \quad (2.453)$$

which easily reduces to

$$I_1 V_1 = I_2 V_2. \quad (2.454)$$

Equations (2.447), (2.450), and (2.454) yield

$$I_1 V_1 = \omega (L_1 I_1^2 + M I_1 I_2) = \omega (L_2 I_2^2 + M I_1 I_2) = I_2 V_2, \quad (2.455)$$

which gives

$$\omega L_1 I_1^2 = \omega L_2 I_2^2, \quad (2.456)$$

and, hence,

$$\frac{I_1}{I_2} = \sqrt{\frac{L_2}{L_1}}. \quad (2.457)$$

Equations (2.454) and (2.457) can be combined to give

$$\frac{V_1}{V_2} = \sqrt{\frac{L_1}{L_2}}. \quad (2.458)$$

Note that, although the mutual inductance of the two coils is entirely responsible for the transfer of energy between the primary and secondary circuits, it is the self inductances of the two coils that determine the ratio of the peak voltages and peak currents in these circuits.

Now, from Section 2.3.5, the self inductances of the primary and secondary coils are given by $L_1 = \mu_0 N_1^2 A/l$ and $L_2 = \mu_0 N_2^2 A/l$, respectively. It follows that

$$\frac{L_1}{L_2} = \left(\frac{N_1}{N_2} \right)^2, \quad (2.459)$$

and, hence, that

$$\frac{V_1}{V_2} = \frac{I_2}{I_1} = \frac{N_1}{N_2}. \quad (2.460)$$

In other words, the ratio of the peak voltages and peak currents in the primary and secondary circuits is determined by the ratio of the number of turns in the primary and secondary coils. This latter ratio is usually called the *turns-ratio* of the transformer. If the secondary coil contains more turns than the primary coil then the peak voltage in the secondary circuit exceeds that in the

primary circuit. This type of transformer is called a *step-up transformer*, because it steps up the voltage of an AC signal. Note that, in a step-up transformer, the peak current in the secondary circuit is less than the peak current in the primary circuit (as must be the case if energy is to be conserved). Thus, a step-up transformer actually steps down the current. Likewise, if the secondary coil contains fewer turns than the primary coil then the peak voltage in the secondary circuit is less than that in the primary circuit. This type of transformer is called a *step-down transformer*. Note that a step-down transformer actually steps up the current (i.e., the peak current in the secondary circuit exceeds that in the primary circuit).

AC electricity is generated in power stations at a fairly low peak voltage (i.e., something like 440V), and is consumed by the domestic user at a peak voltage of 110V (in the U.S.). However, AC electricity is transmitted from the power station to the location where it is consumed at a very high peak voltage (typically 50kV). In fact, as soon as an AC signal comes out of a generator in a power station it is fed into a step-up transformer that boosts its peak voltage from a few hundred volts to many tens of kilovolts. The output from the step-up transformer is fed into a high tension transmission line, which typically transports the electricity over many tens of kilometers, and, once the electricity has reached its point of consumption, it is fed through a series of step-down transformers, until, by the time it emerges from a domestic power socket, its peak voltage is only 110V. But, if AC electricity is both generated and consumed at comparatively low peak voltages, why go to the trouble of stepping up the peak voltage to a very high value at the power station, and then stepping down the voltage again once the electricity has reached its point of consumption? Why not generate, transmit, and distribute the electricity at a peak voltage of 110V? Well, consider an electric power line that transmits a peak electric power P between a power station and a city. We can think of P , which depends on the number of consumers in the city, and the nature of the electrical devices that they operate, as essentially a fixed number. Suppose that V and I are the peak voltage and peak current of the AC signal transmitted along the transmission line, respectively. We can think of these numbers as being variable, because we can change them using a transformer. However, because $P = IV$, the product of the peak voltage and the peak current must remain constant. Suppose that the resistance of the transmission line is R . The peak rate at which electrical energy is lost due to ohmic heating in the line is $P_R = I^2 R$ (see Section 2.1.11), which can be written

$$P_R = \frac{P^2 R}{V^2}. \quad (2.461)$$

Thus, if the power, P , transmitted down the line is a fixed quantity, as is the resistance, R , of the line, then the power lost in the line due to ohmic heating varies like the inverse square of the peak voltage in the line. It turns out that even at very high voltages, such as 50kV, the ohmic power losses in transmission lines that run over tens of kilometers can amount to up to 20% of the transmitted power. It can readily be appreciated that if an attempt were made to transmit AC electric power at a peak voltage of 110V then the ohmic losses would be so severe that virtually none of the power would reach its destination. Thus, it is only possible to generate electric power at a central location, transmit it over long distances, and then distribute it at its point of consumption, if the transmission is performed at a very high peak voltages (in fact, the higher, the better). Transformers play a vital role in this process because they allow us to step-up and step-down the voltage of an AC electric signal very efficiently (a well-designed transformer typically has a power loss that is only a few

percent of the total power flowing through it).

Of course, transformers do not work for direct current (DC) electricity, because the magnetic flux generated by the primary coil does not vary in time, and, therefore, does not induce an emf in the secondary coil. In fact, there is no efficient method of stepping-up or stepping-down the voltage of a DC electric signal. Thus, it is impossible to efficiently transmit DC electric power over long distances. This is the main reason why commercially generated electricity is AC, rather than DC.

2.4 Maxwell's Equations

2.4.1 Displacement Current

Michael Faraday revolutionized physics in 1830 by showing that electricity and magnetism were interrelated phenomena. (See Section 2.3.1.) He achieved this breakthrough by careful experimentation. Between 1864 and 1873, James Clerk Maxwell achieved a similar breakthrough by pure thought. Of course, this was only possible because he was able to take the previous experimental results of Coulomb, Ampère, Faraday, et cetera, as his starting point.

Prior to 1864, the laws of electromagnetism were written in integral form. Thus, Gauss's law (in SI units) was expressed as follows; the flux of the electric field, $\mathbf{E}(\mathbf{r}, t)$, through a closed surface, S , enclosing a volume, V , is equal to the net enclosed electric charge, divided by ϵ_0 ; or

$$\oint_S \mathbf{E} \cdot d\mathbf{S} = \frac{1}{\epsilon_0} \int_V \rho(\mathbf{r}, t) dV, \quad (2.462)$$

where $\rho(\mathbf{r}, t)$ is the electric charge density. (See Section 2.1.6.) The no magnetic monopole law was expressed as follows; the flux of the magnetic field, $\mathbf{B}(\mathbf{r}, t)$, through any closed surface, S is zero; or

$$\oint_S \mathbf{B} \cdot d\mathbf{S} = 0. \quad (2.463)$$

(See Section 2.2.9.) Faraday's law of electromagnetic induction was expressed as follows; the line integral of the electric field around a closed loop, C , is equal to minus the rate of change of the magnetic flux passing through any surface, S , attached to the loop; or

$$\oint_C \mathbf{E} \cdot d\mathbf{r} = -\frac{\partial}{\partial t} \int_S \mathbf{B} \cdot d\mathbf{S}. \quad (2.464)$$

(See Section 2.3.1.) Finally, Ampère's circuital law was expressed as follows; the line integral of the magnetic field around a closed loop C is equal the net current passing through any surface, S , attached to the loop, multiplied by μ_0 ; or

$$\oint_C \mathbf{B} \cdot d\mathbf{r} = \mu_0 \int_S \mathbf{j}(\mathbf{r}, t) \cdot d\mathbf{S}, \quad (2.465)$$

where $\mathbf{j}(\mathbf{r}, t)$ is the electric current density. (See Section 2.2.10.)

Maxwell's first great achievement was to realize that, with the aid of the divergence theorem and the curl theorem (see Sections A.20 and A.22), these laws could be re-expressed as a set of

first-order partial differential equations. Of course, he wrote his equations out in component form, because modern vector notation did not come into vogue until about the time of the First World War. In modern notation, Maxwell first wrote:

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}, \quad (2.466)$$

$$\nabla \cdot \mathbf{B} = 0, \quad (2.467)$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}, \quad (2.468)$$

$$\nabla \times \mathbf{B} = \mu_0 \mathbf{j}. \quad (2.469)$$

[See Equations (2.54), (2.263), (2.286), and (2.271).] Maxwell's second great achievement was to realize that these equations are not mathematically self-consistent.

Consider the integral form of Equation (2.469):

$$\oint_C \mathbf{B} \cdot d\mathbf{r} = \mu_0 \int_S \mathbf{j} \cdot d\mathbf{S}. \quad (2.470)$$

This equation states that the line integral of the magnetic field around a closed loop C is equal to the flux of the current density through the loop, multiplied by μ_0 . The problem is that the flux of the current density through a loop is not, in general, a well-defined quantity. In order for the flux to be well defined, the integral of $\mathbf{j} \cdot d\mathbf{S}$ over some surface S attached to a loop C must depend on C , but not on the details of S . This is only the case if

$$\nabla \cdot \mathbf{j} = 0. \quad (2.471)$$

(See Section A.20.) Unfortunately, the previous condition is only satisfied for non-time-varying fields.

Why do we say that, in general, $\nabla \cdot \mathbf{j} \neq 0$? Consider the flux of \mathbf{j} out of some closed surface, S , enclosing a volume, V . This is clearly equivalent to the instantaneous rate at which electric charge flows out of S . However, because electric charge is a conserved quantity (see Section 2.1.2), the rate at which charge flows out of S must equal the rate of decrease of the charge contained in volume V . Thus,

$$\oint_S \mathbf{j} \cdot d\mathbf{S} = -\frac{\partial}{\partial t} \int_V \rho dV. \quad (2.472)$$

Making use of the divergence theorem (see Section A.20), the previous equation yields

$$\nabla \cdot \mathbf{j} = -\frac{\partial \rho}{\partial t}. \quad (2.473)$$

Thus, it is only the case that $\nabla \cdot \mathbf{j} = 0$ in a steady state situation; that is, when $\partial/\partial t \equiv 0$.

The problem with Ampère's circuital law is well illustrated by the following very famous example. Consider a long straight wire interrupted by a parallel plate capacitor. Suppose that C is some loop that circles the wire. In the time-independent case, the capacitor acts like a break in the

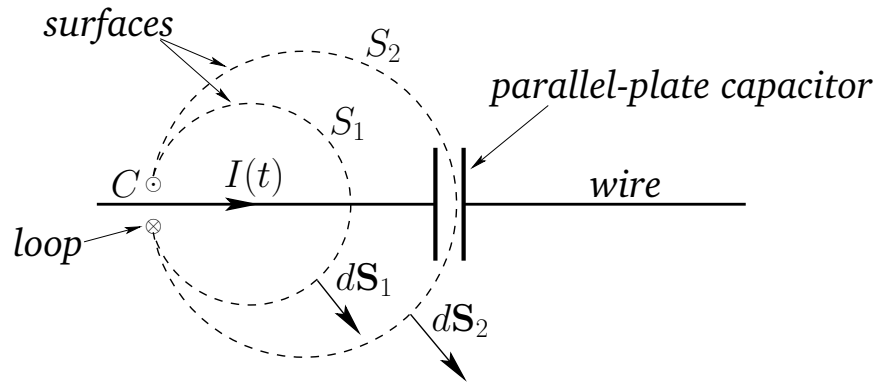


Figure 2.39: Application of Ampère's circuital law to a charging, or discharging, capacitor.

wire, so no current flows, and no magnetic field is generated. There is clearly no problem with Ampère's circuital law in this case. However, in the time-dependent case, a transient current flows in the wire as the capacitor charges up, or charges down, and so a transient magnetic field is generated. Thus, the line integral of the magnetic field around C is (transiently) non-zero. According to Ampère's circuital law, the flux of the current density through any surface attached to C should also be (transiently) non-zero. Let us consider two such surfaces. The first surface, S_1 , intersects the wire. See Figure 2.39. This surface causes us no problem, because the flux of \mathbf{j} through the surface is clearly non-zero (because the surface intersects a current-carrying wire). The second surface, S_2 , passes between the plates of the capacitor, and, therefore, does not intersect the wire at all. Clearly, the flux of the current density through this surface is zero. The current density fluxes through surfaces S_1 and S_2 are obviously different. However, both surfaces are attached to the same loop C , so the fluxes should be the same, according to Ampère's circuital law, (2.470). Note, however, that although the surface S_2 does not intersect any electric current, it does pass through a region containing a strong, time-varying electric field, as it threads between the plates of the charging (or discharging) capacitor. Perhaps, if we add a term involving $\partial\mathbf{E}/\partial t$ to the right-hand side of Equation (2.469) then we can somehow fix up Ampère's circuital law? This is, essentially, how Maxwell reasoned one hundred and fifty years ago.

Let us try out this scheme. Suppose that we write

$$\nabla \times \mathbf{B} = \mu_0 \mathbf{j} + \lambda \frac{\partial \mathbf{E}}{\partial t}, \quad (2.474)$$

instead of Equation (2.469). Here, λ is some constant. Does this resolve our problem? We require the flux of the right-hand side of the previous equation through some loop C to be well defined; that is, the flux should only depend on C , and not the particular surface S (which spans C) upon which it is evaluated. This is another way of saying that we require the divergence of the right-hand side of the previous equation to be zero. (See Section A.20.) In fact, we can see that this is necessary for mathematical self-consistency, because the divergence of the left-hand side is identically zero. (See Section A.22.) So, taking the divergence of Equation (2.474), we obtain

$$0 = \mu_0 \nabla \cdot \mathbf{j} + \lambda \frac{\partial \nabla \cdot \mathbf{E}}{\partial t}. \quad (2.475)$$

But, we know that

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0} \quad (2.476)$$

[see Equation (2.466)], so combining the previous two equations we arrive at

$$\mu_0 \nabla \cdot \mathbf{j} + \frac{\lambda}{\epsilon_0} \frac{\partial \rho}{\partial t} = 0. \quad (2.477)$$

Now, our charge conservation law, (2.473), can be written

$$\nabla \cdot \mathbf{j} + \frac{\partial \rho}{\partial t} = 0. \quad (2.478)$$

The previous two equations are in agreement provided $\lambda = \epsilon_0 \mu_0$. So, if we modify Equation (2.469) such that it reads

$$\nabla \times \mathbf{B} = \mu_0 (\mathbf{j} + \mathbf{j}_d), \quad (2.479)$$

where

$$\mathbf{j}_d = \epsilon_0 \frac{\partial \mathbf{E}}{\partial t}, \quad (2.480)$$

then we find that the divergence of the right-hand side is zero, as a consequence of charge conservation. The additional term, \mathbf{j}_d , is known as the *displacement current* density (this name was invented by Maxwell). In summary, we have shown that, although the flux of the real current density through a loop is not well defined, if we form the sum of the real current density and the displacement current density then the flux of this new quantity through a loop is well defined.

Of course, the displacement current is not a current at all. It is, in fact, associated with the induction of magnetic fields by time-varying electric fields. Maxwell came up with this rather curious name because many of his ideas regarding electric and magnetic fields were completely wrong. For instance, Maxwell believed in the aether (a tenuous invisible medium permeating all space; see Section 3.1.2), and he thought that electric and magnetic fields corresponded to stresses in this medium. He also thought that the displacement current was associated with a displacement of the aether (hence, the name). The reason that these misconceptions did not invalidate Maxwell's equations is quite simple. Maxwell based his equations on the results of experiments, and he added in his extra term so as to make these equations mathematically self-consistent. Both of these steps are valid irrespective of the existence or non-existence of the aether.

The field equations (2.466)–(2.469) are derived directly from the results of famous nineteenth century experiments. So, if a new term involving the time derivative of the electric field needs to be added to one of these equations, for the sake of mathematical consistency, why is there is no corresponding nineteenth century experimental result that demonstrates this fact? Actually, as is described in the following, the new term corresponds to an effect that is far too small to have been observed in the nineteenth century.

First, we shall show that it is comparatively easy to detect the induction of an electric field by a changing magnetic field in a desktop laboratory experiment. The Earth's magnetic field is about 1 gauss (that is, 10^{-4} tesla). Magnetic fields generated by electromagnets (that will fit on a laboratory desktop) are typically about one hundred times larger than this. Let us, therefore,

consider a hypothetical experiment in which a 100 gauss magnetic field is switched on suddenly. Suppose that the field ramps up in one tenth of a second. What electromotive force is generated in a 10 centimeter square loop of wire located in this field? Faraday's law is written

$$V = -\frac{\partial}{\partial t} \oint \mathbf{B} \cdot d\mathbf{S} \approx \frac{BA}{t}, \quad (2.481)$$

where $B = 0.01$ tesla is the magnetic field-strength, $A = 0.01 \text{ m}^2$ the area of the loop, and $t = 0.1$ seconds the ramp time. (See Section 2.3.1.) It follows that $V \approx 1$ millivolt, which is easily detectable. In fact, most hand-held laboratory voltmeters are calibrated in millivolts. It is, thus, clear that we would have no difficulty whatsoever detecting the magnetic induction of electric fields in a nineteenth-century-style laboratory experiment.

Let us now consider the electric induction of magnetic fields. Suppose that our electric field is generated by a parallel plate capacitor of spacing one centimeter that is charged up to 100 volts. This gives an electric field of 10^4 volts per meter. Suppose, further, that the capacitor is discharged in one tenth of a second. The law of electric induction is obtained by integrating Equation (2.479), and neglecting the first term on the right-hand side. Thus,

$$\oint_C \mathbf{B} \cdot d\mathbf{r} = \epsilon_0 \mu_0 \frac{\partial}{\partial t} \int_S \mathbf{E} \cdot d\mathbf{S}. \quad (2.482)$$

Let us consider a loop that is 10 centimeters square. What is the magnetic field generated around this loop (which we could try to measure with a Hall probe)? Very approximately, we find that

$$lB \approx \epsilon_0 \mu_0 \frac{El^2}{t}, \quad (2.483)$$

where $l = 0.1$ meters is the dimensions of the loop, B the magnetic field-strength, $E = 10^4$ volts per meter the electric field, and $t = 0.1$ seconds the decay time of the field. We obtain $B \approx 10^{-9}$ gauss. Modern technology is unable to detect such a small magnetic field, so we cannot really blame nineteenth century physicists for not discovering electric induction experimentally.

Note, however, that the displacement current is detectable in some modern experiments. Suppose that we take an FM radio signal, amplify it so that its peak voltage is one hundred volts, and then apply it to the parallel plate capacitor in the previous hypothetical experiment. What size of magnetic field would this generate? A typical FM signal oscillates at 10^9 Hz, so t in the previous example changes from 0.1 seconds to 10^{-9} seconds. Thus, the induced magnetic field is about 10^{-1} gauss. This is certainly detectable by modern technology. Hence, we conclude that if the electric field is oscillating sufficiently rapidly then electric induction of magnetic fields is an observable effect. In fact, there is a virtually infallible rule for deciding whether or not the displacement current can be neglected in Equation (2.479). Namely, if electromagnetic radiation is important then the displacement current must be included. On the other hand, if electromagnetic radiation is unimportant then the displacement current can be safely neglected. Clearly, Maxwell's inclusion of the displacement current in Equation (2.479) was a vital step in his later realization that his equations allowed propagating wave-like solutions. These solutions are, of course, electromagnetic waves.

2.4.2 Maxwell's Equations

We are now in a position to write out a complete, mathematically self-consistent, set of field equations that govern electric and magnetic phenomena:

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}, \quad (2.484)$$

$$\nabla \cdot \mathbf{B} = 0, \quad (2.485)$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}, \quad (2.486)$$

$$\nabla \times \mathbf{B} = \mu_0 \mathbf{j} + \epsilon_0 \mu_0 \frac{\partial \mathbf{E}}{\partial t}. \quad (2.487)$$

The equations are known as *Maxwell's equation*. The first Maxwell equation describes how electric fields are induced by electric charges, and is equivalent to Gauss's law. The second Maxwell equation states that there is no such thing as a magnetic monopole. The third Maxwell equation describes the induction of electric fields by changing magnetic fields, and is equivalent to Faraday's law of electromagnetic induction. The fourth Maxwell equation describes the generation of magnetic fields by electric currents, and the induction of magnetic fields by changing electric fields, and incorporates Ampère's circuital law.

As an example of a calculation involving the displacement current, let us find the current and displacement current densities associated with the decaying charge distribution

$$\rho(r, t) = \frac{\rho_0 \exp(-t/\tau)}{r^2 + a^2}, \quad (2.488)$$

where r is a spherical polar coordinate (see Section A.23), ρ_0 is a constant, and τ and a are positive constants. Now, according to charge conservation,

$$\nabla \cdot \mathbf{j} = -\frac{\partial \rho}{\partial t}. \quad (2.489)$$

[See Equation (2.473).] By symmetry, we expect $\mathbf{j} = \mathbf{j}(r, t)$. Hence, it follows that $\mathbf{j} = j_r(r, t) \mathbf{e}_r$ [because only a radial current has a non-zero divergence when $\mathbf{j} = \mathbf{j}(r)$]. [See Equation (A.173).] Thus, the previous two equations yield.

$$\frac{1}{r^2} \frac{\partial(r^2 j_r)}{\partial r} = -\frac{\partial \rho}{\partial t} = \frac{\rho_0 \exp(-t/\tau)}{\tau(r^2 + a^2)}. \quad (2.490)$$

The previous expression can be integrated, subject to the sensible boundary condition $j_r(0) = 0$, to give

$$j_r(r) = \frac{\rho_0}{\tau} e^{-t/\tau} \left[\frac{r - a \tan^{-1}(r/a)}{r^2} \right]. \quad (2.491)$$

Now, the electric field generated by the decaying charge distribution satisfies

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}. \quad (2.492)$$

[See Equation (2.466).] Because $\partial\rho/\partial t = -\rho/\tau$, it can be seen, from a comparison of Equations (2.489) and (2.492), that

$$\mathbf{E} = \frac{\tau}{\epsilon_0} \mathbf{j}. \quad (2.493)$$

However, the displacement current density is given by

$$\mathbf{j}_d = \epsilon_0 \frac{\partial \mathbf{E}}{\partial t} = -\mathbf{j}, \quad (2.494)$$

because $\partial \mathbf{j} / \partial t = -\mathbf{j} / \tau$. Hence, we conclude that the displacement current density cancels out the true current density, so that $\mathbf{j} + \mathbf{j}_d = \mathbf{0}$. This is just as well, because $\nabla \times \mathbf{B} = \mu_0 (\mathbf{j} + \mathbf{j}_d)$. [See Equation (2.479).] But, if $\mathbf{B} = \mathbf{B}(r, t)$ then, by symmetry, $\nabla \times \mathbf{B}$ has no radial component. [See Equation (A.174).] Thus, if the current and displacement current are constrained, by symmetry, to be radial, then they must sum to zero, otherwise the fourth Maxwell equation cannot be satisfied. In fact, no magnetic field is generated in this particular example.

2.4.3 Potential Formulation of Maxwell's Equations

We saw in Section 2.3.2 that the second and third Maxwell equations, (2.485) and (2.486), are automatically satisfied if we write the electric and magnetic fields in terms of scalar and vector potentials; that is,

$$\mathbf{E} = -\nabla\phi - \frac{\partial \mathbf{A}}{\partial t}, \quad (2.495)$$

$$\mathbf{B} = \nabla \times \mathbf{A}. \quad (2.496)$$

As was discussed in Section 2.3.3, this prescription is not unique, but we can make it unique by adopting the following conventions:

$$\phi(\mathbf{r}) \rightarrow 0 \quad \text{as } |\mathbf{r}| \rightarrow \infty, \quad (2.497)$$

$$\nabla \cdot \mathbf{A} = -\epsilon_0 \mu_0 \frac{\partial \phi}{\partial t}. \quad (2.498)$$

The previous equation is known as the *Lorenz gauge*.

The previous equation can be combined with Equation (2.495) and the first Maxwell equation, (2.484), to give

$$\epsilon_0 \mu_0 \frac{\partial^2 \phi}{\partial t^2} - \nabla^2 \phi = \frac{\rho}{\epsilon_0}. \quad (2.499)$$

Let us now consider the fourth Maxwell equation, (2.487). Substitution of Equations (2.495) and (2.496) into this equation yields

$$\nabla \times \nabla \times \mathbf{A} \equiv \nabla(\nabla \cdot \mathbf{A}) - \nabla^2 \mathbf{A} = \mu_0 \mathbf{j} - \epsilon_0 \mu_0 \frac{\partial \nabla \phi}{\partial t} - \epsilon_0 \mu_0 \frac{\partial^2 \mathbf{A}}{\partial t^2} \quad (2.500)$$

(see Section A.24), or

$$\epsilon_0 \mu_0 \frac{\partial^2 \mathbf{A}}{\partial t^2} - \nabla^2 \mathbf{A} = \mu_0 \mathbf{j} - \nabla \left(\nabla \cdot \mathbf{A} + \epsilon_0 \mu_0 \frac{\partial \phi}{\partial t} \right). \quad (2.501)$$

We can now see quite clearly where the Lorenz gauge, (2.498), comes from. The previous equation is, in general, very complicated, because it involves both the vector and scalar potentials. However, if we adopt the Lorenz gauge then the last term on the right-hand side becomes zero, and the equation simplifies considerably, and ends up only involving the vector potential. Thus, we find that Maxwell's equations reduce to the following equations:

$$\epsilon_0 \mu_0 \frac{\partial^2 \phi}{\partial t^2} - \nabla^2 \phi = \frac{\rho}{\epsilon_0}, \quad (2.502)$$

$$\epsilon_0 \mu_0 \frac{\partial^2 \mathbf{A}}{\partial t^2} - \nabla^2 \mathbf{A} = \mu_0 \mathbf{j}. \quad (2.503)$$

Of course, this is the same (scalar) equation written four times over. In a non-time-varying situation (i.e., $\partial/\partial t = 0$), the equation in question reduces to Poisson's equation (see Section 2.1.9), which we know how to solve. With the $\partial^2/\partial t^2$ terms included, the equation becomes a slightly more complicated equation (in fact, it is an inhomogeneous three-dimensional wave equation).

2.4.4 Electromagnetic Waves

Let us demonstrate that Maxwell's equations possess wave-like solutions that can propagate through a vacuum. These solutions are known as *electromagnetic waves*. Let us start from Maxwell's equations in free space (i.e., with no charges and no currents):

$$\nabla \cdot \mathbf{E} = 0, \quad (2.504)$$

$$\nabla \cdot \mathbf{B} = 0, \quad (2.505)$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}, \quad (2.506)$$

$$\nabla \times \mathbf{B} = \epsilon_0 \mu_0 \frac{\partial \mathbf{E}}{\partial t}. \quad (2.507)$$

[See Equations (2.484)–(2.487).]

There is an easy way to show that the previous equations possess wave-like solutions, and a hard way. The easy way is to assume that the solutions are going to be wave-like beforehand. Specifically, let us search for plane-wave solutions of the form:

$$\mathbf{E}(\mathbf{r}, t) = \mathbf{E}_0 \cos(\mathbf{k} \cdot \mathbf{r} - \omega t), \quad (2.508)$$

$$\mathbf{B}(\mathbf{r}, t) = \mathbf{B}_0 \cos(\mathbf{k} \cdot \mathbf{r} - \omega t - \phi). \quad (2.509)$$

Here, \mathbf{E}_0 and \mathbf{B}_0 are constant vectors, \mathbf{k} is known as the *wavevector*, and ω is the angular frequency of oscillation of the wave. The frequency in hertz, f , is related to the angular frequency via

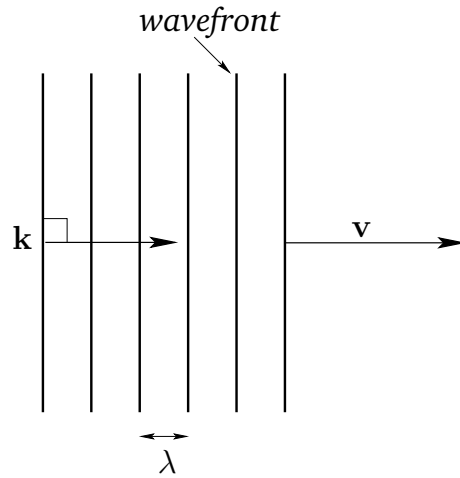


Figure 2.40: Wavefronts associated with a plane wave.

$\omega = 2\pi f$; this frequency is conventionally defined to be positive. The quantity ϕ is a phase difference between the electric and magnetic fields. Actually, it is more convenient to write

$$\mathbf{E}(\mathbf{r}, t) = \mathbf{E}_0 e^{i(\mathbf{k} \cdot \mathbf{r} - \omega t)}, \quad (2.510)$$

$$\mathbf{B}(\mathbf{r}, t) = \mathbf{B}_0 e^{i(\mathbf{k} \cdot \mathbf{r} - \omega t)}, \quad (2.511)$$

where, by convention, the physical solution is the real part of the previous equations. The phase difference ϕ is absorbed into the constant vector \mathbf{B}_0 by allowing it to become complex. Thus, $\mathbf{B}_0 \rightarrow \mathbf{B}_0 e^{-i\phi}$. In general, the vector \mathbf{E}_0 is also complex.

Now, assuming (without loss of generality) that \mathbf{E}_0 is real, a wave maximum of the electric field satisfies

$$\mathbf{k} \cdot \mathbf{r} = \omega t + n2\pi, \quad (2.512)$$

where n is an integer. The solution to this equation is a set of equally-spaced parallel planes (one plane for each possible value of n), whose normals are parallel to the wavevector \mathbf{k} , and that propagate in the direction of \mathbf{k} with phase velocity

$$c = \frac{\omega}{k}. \quad (2.513)$$

The spacing between adjacent planes (i.e., the wavelength) is given by

$$\lambda = \frac{2\pi}{k}. \quad (2.514)$$

See Figure 2.40.

Consider a general plane-wave vector field

$$\mathbf{A}(\mathbf{r}, t) = \mathbf{A}_0 e^{i(\mathbf{k} \cdot \mathbf{r} - \omega t)}. \quad (2.515)$$

What is the divergence of \mathbf{A} ? This is easy to evaluate. We have

$$\nabla \cdot \mathbf{A} \equiv \frac{\partial A_x}{\partial x} + \frac{\partial A_y}{\partial y} + \frac{\partial A_z}{\partial z} = (A_{0x} i k_x + A_{0y} i k_y + A_{0z} i k_z) e^{i(\mathbf{k} \cdot \mathbf{r} - \omega t)} = i \mathbf{k} \cdot \mathbf{A}. \quad (2.516)$$

(See Section A.20.) How about the curl of \mathbf{A} ? This is slightly more difficult. We have

$$(\nabla \times \mathbf{A})_x \equiv \frac{\partial A_z}{\partial y} - \frac{\partial A_y}{\partial z} = (i k_y A_z - i k_z A_y) = i (\mathbf{k} \times \mathbf{A})_x \quad (2.517)$$

(see Section A.22), which easily generalizes to

$$\nabla \times \mathbf{A} = i \mathbf{k} \times \mathbf{A}. \quad (2.518)$$

Hence, it is apparent that vector field operations on a plane-wave vector field are equivalent to replacing the ∇ operator with $i \mathbf{k}$. Of course, the $\partial/\partial t$ operator can be replaced by $-i \omega$.

The first Maxwell equation, (2.504), reduces to

$$i \mathbf{k} \cdot \mathbf{E}_0 = 0, \quad (2.519)$$

using the assumed electric and magnetic fields, (2.510) and (2.511), and Equation (2.516). Thus, the electric field is perpendicular to the direction of propagation of the wave. (See Section A.6.) Likewise, the second Maxwell equation, (2.505), gives

$$i \mathbf{k} \cdot \mathbf{B}_0 = 0, \quad (2.520)$$

implying that the magnetic field is also perpendicular to the direction of propagation. Clearly, the wave-like solutions of Maxwell's equation are a type of transverse wave. The third Maxwell equation, (2.506), yields

$$i \mathbf{k} \times \mathbf{E}_0 = i \omega \mathbf{B}_0, \quad (2.521)$$

where use has been made of Equation (2.518). Forming the scalar product of this equation with \mathbf{E}_0 gives

$$\mathbf{E}_0 \cdot \mathbf{B}_0 = \frac{\mathbf{E}_0 \cdot \mathbf{k} \times \mathbf{E}_0}{\omega} = 0. \quad (2.522)$$

Thus, the electric and magnetic fields are mutually perpendicular. (See Sections A.6 and A.10.) Forming the scalar product of Equation (2.521) with \mathbf{B}_0 yields

$$\mathbf{B}_0 \cdot \mathbf{k} \times \mathbf{E}_0 = \omega B_0^2 > 0. \quad (2.523)$$

Thus, the vectors \mathbf{E}_0 , \mathbf{B}_0 , and \mathbf{k} are mutually perpendicular, and form a right-handed set. (See Section A.10.) The final Maxwell equation, (2.507), gives

$$i \mathbf{k} \times \mathbf{B}_0 = -i \epsilon_0 \mu_0 \omega \mathbf{E}_0. \quad (2.524)$$

Combining this equation with Equation (2.521) yields

$$\mathbf{k} \times (\mathbf{k} \times \mathbf{E}_0) \equiv (\mathbf{k} \cdot \mathbf{E}_0) \mathbf{k} - k^2 \mathbf{E}_0 = -k^2 \mathbf{E}_0 = -\epsilon_0 \mu_0 \omega^2 \mathbf{E}_0, \quad (2.525)$$

or

$$k^2 = \epsilon_0 \mu_0 \omega^2, \quad (2.526)$$

where use has been made of Equation (2.519). (See Section A.11.) However, we know, from Equation (2.513), that the phase velocity, c , of the wave is related to the magnitude of the wavevector and the angular wave frequency via $c = \omega/k$. Thus, we obtain

$$c = \frac{1}{\sqrt{\epsilon_0 \mu_0}}. \quad (2.527)$$

We have found transverse plane-wave solutions of the free-space Maxwell equations propagating at some phase velocity c , that is given by a combination of ϵ_0 and μ_0 , and is, thus, the same for all frequencies and wavelengths. The constants ϵ_0 and μ_0 are easily measurable. The former is related to the force acting between stationary electric charges, and the latter to the force acting between steady electric currents. Both of these constants were fairly well known in Maxwell's time. Maxwell, incidentally, was the first person to look for wave-like solutions of his equations, and, thus, to derive Equation (2.527). The modern values of ϵ_0 and μ_0 are

$$\epsilon_0 = 8.8542 \times 10^{-12} \text{ C}^2 \text{ N}^{-1} \text{ m}^{-2}, \quad (2.528)$$

$$\mu_0 = 4\pi \times 10^{-7} \text{ N A}^{-2}. \quad (2.529)$$

Let us use these values to find the phase velocity of electromagnetic waves. We obtain

$$c = \frac{1}{\sqrt{\epsilon_0 \mu_0}} = 2.998 \times 10^8 \text{ m s}^{-1}. \quad (2.530)$$

Of course, we immediately recognize this as the speed of light in vacuum. Maxwell also made this connection back in the 1870's. He conjectured that light, whose nature had previously been unknown, was a form of electromagnetic radiation. This was a remarkable prediction. After all, Maxwell's equations were derived from the results of bench-top laboratory experiments involving charges, batteries, coils, and currents, that apparently had nothing whatsoever to do with light.

Maxwell was able to make another remarkable prediction. The wavelength of light was well-known in the late nineteenth century from studies of diffraction through slits, et cetera. Visible light actually occupies a surprisingly narrow wavelength range. The shortest wavelength blue light that is visible to the typical human eye has a wavelength of $\lambda = 0.38$ microns (one micron is 10^{-6} meters). The longest wavelength red light that is visible has a wavelength of $\lambda = 0.75$ microns. However, there is nothing in our analysis that suggests that this particular range of wavelengths is special. Electromagnetic waves can have any wavelength. Maxwell concluded that visible light was a small part of a vast spectrum of previously undiscovered types of electromagnetic radiation. Since Maxwell's time, virtually all of the non-visible parts of the electromagnetic spectrum have been observed.

Table 1 gives a brief guide to the electromagnetic spectrum. Electromagnetic waves are of particular importance to us because they are our main source of information regarding the universe around us. Radio waves and microwaves (which are comparatively hard to scatter) have provided much of our knowledge about the center of our own galaxy. This is completely unobservable in

Radiation type	Wavelength range (m)
Gamma Rays	$< 10^{-11}$
X-Rays	$10^{-11} - 10^{-9}$
Ultraviolet	$10^{-9} - 10^{-7}$
Visible	$10^{-7} - 10^{-6}$
Infrared	$10^{-6} - 10^{-4}$
Microwave	$10^{-4} - 10^{-1}$
TV-FM	$10^{-1} - 10^1$
Radio	$> 10^1$

Table 2.1: The electromagnetic spectrum

visible light, which is strongly scattered by interstellar gas and dust lying in the galactic plane. For the same reason, the spiral arms of our galaxy can only be mapped out using radio waves. Infrared radiation is useful for detecting protostars, which are not yet hot enough to emit visible radiation. Of course, visible radiation is still the mainstay of astronomy. Satellite-based ultraviolet observations have yielded invaluable insights into the structure and distribution of distant galaxies. Finally, X-ray and γ -ray astronomy usually concentrates on exotic objects, such as pulsars and supernova remnants.

Equations (2.519), (2.521), and the relation $c = \omega/k$, imply that

$$B_0 = \frac{E_0}{c}. \quad (2.531)$$

Thus, the magnetic field associated with an electromagnetic wave is smaller in magnitude than the electric field by a factor c . Consider an electrically charged particle interacting with an electromagnetic wave. The force exerted on the particle is given by the Lorentz force law,

$$\mathbf{f} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}). \quad (2.532)$$

(See Section 2.2.4.) The ratio of the electric and magnetic forces is

$$\frac{f_{\text{magnetic}}}{f_{\text{electric}}} \simeq \frac{v B_0}{E_0} \simeq \frac{v}{c}. \quad (2.533)$$

So, unless the particle is moving close to the speed of light (i.e., unless the particle is relativistic), the electric force greatly exceeds the magnetic force. Clearly, in most terrestrial situations, electromagnetic waves are an essentially electrical phenomenon (as far as their interaction with matter is concerned). For this reason, electromagnetic waves are usually characterized by their wavevector, \mathbf{k} (which specifies the direction of propagation and the wavelength), and the plane of *polarization* (i.e., the plane of oscillation) of the associated electric field. For a given wavevector, \mathbf{k} , the electric field can have any direction in the plane normal to \mathbf{k} . [See Equation (2.519).] However, there are only two independent directions in a plane (i.e., we can only define two linearly

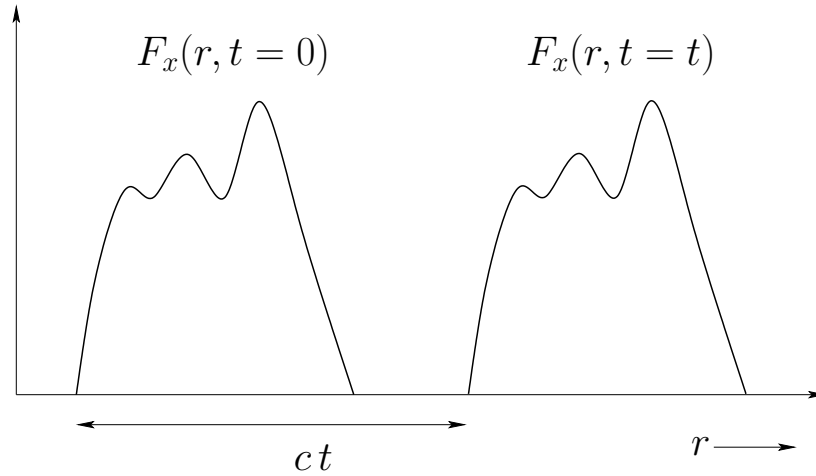


Figure 2.41: An arbitrary wave-pulse.

independent vectors in a plane). This implies that there are only two independent polarizations of an electromagnetic wave, once its direction of propagation is specified.

But, how do electromagnetic waves propagate through a vacuum? After all, most types of wave require a medium before they can propagate (e.g., sound waves require air). The answer to this question is evident from Equations (2.506) and (2.507). According to these equations, the time variation of the electric component of the wave induces the magnetic component, and the time variation of the magnetic component induces the electric component. In other words, electromagnetic waves are self-sustaining, and, therefore, require no medium through which to propagate.

Let us now search for the wave-like solutions of Maxwell's equations in free-space the hard way. Suppose that we take the curl of the fourth Maxwell equation, (2.507). We obtain

$$\nabla \times \nabla \times \mathbf{B} \equiv \nabla(\nabla \cdot \mathbf{B}) - \nabla^2 \mathbf{B} = -\nabla^2 \mathbf{B} = \epsilon_0 \mu_0 \frac{\partial \nabla \times \mathbf{E}}{\partial t}. \quad (2.534)$$

[See Equation (A.187).] Here, we have made use of the fact that $\nabla \cdot \mathbf{B} = 0$, according to the second Maxwell equation, (2.505). The third Maxwell equation, (2.506), yields

$$\left(\frac{1}{c^2} \frac{\partial^2}{\partial t^2} - \nabla^2 \right) \mathbf{B} = \mathbf{0}, \quad (2.535)$$

where use has been made of Equation (2.530). A similar equation can be obtained for the electric field by taking the curl of Equation (2.506):

$$\left(\frac{1}{c^2} \frac{\partial^2}{\partial t^2} - \nabla^2 \right) \mathbf{E} = \mathbf{0}, \quad (2.536)$$

We have found that electric and magnetic fields both satisfy equations of the form

$$\left(\frac{1}{c^2} \frac{\partial^2}{\partial t^2} - \nabla^2 \right) \mathbf{A} = \mathbf{0} \quad (2.537)$$

in free space. As is easily verified, the most general solution to this equation is

$$A_x = F_x(\mathbf{n} \cdot \mathbf{r} - ct), \quad (2.538)$$

$$A_y = F_y(\mathbf{n} \cdot \mathbf{r} - ct), \quad (2.539)$$

$$A_z = F_z(\mathbf{n} \cdot \mathbf{r} - ct), \quad (2.540)$$

where \mathbf{n} is a unit vector, and $F_x(\phi)$, $F_y(\phi)$, and $F_z(\phi)$ are arbitrary one-dimensional scalar functions. Looking along the direction of \mathbf{n} , so that $\mathbf{n} \cdot \mathbf{r} = r$, we find that

$$A_x = F_x(r - ct), \quad (2.541)$$

$$A_y = F_y(r - ct), \quad (2.542)$$

$$A_z = F_z(r - ct). \quad (2.543)$$

The x -component of this solution is shown schematically in Figure 2.41. The solution clearly propagates along the r -axis, at the speed c , without changing shape. If we look along a direction that is perpendicular to \mathbf{n} then $\mathbf{n} \cdot \mathbf{r} = 0$, and there is no propagation. Thus, the components of \mathbf{A} are arbitrarily-shaped pulses that propagate, without changing shape, along the direction of \mathbf{n} with speed c . These pulses can be related to the sinusoidal plane-wave solutions which we found earlier by Fourier transformation; that is,

$$F_x(r - ct) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \bar{F}_x(k) e^{ik(r-ct)} dk, \quad (2.544)$$

where

$$\bar{F}_x(k) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} F_x(x) e^{-ikx} dx. \quad (2.545)$$

(See Section 4.2.4.) Thus, any arbitrary-shaped pulse propagating in the direction of \mathbf{n} with speed c can be broken down into a superposition of sinusoidal oscillations of different wavevectors, $k\mathbf{n}$, propagating in the same direction with the same speed.

2.4.5 Energy Conservation

We have seen that the energy density of an electric field is given by [see Equation (2.84)]

$$U_E = \frac{\epsilon_0 E^2}{2}, \quad (2.546)$$

whereas the energy density of a magnetic field takes the form [see Equation (2.377)]

$$U_B = \frac{B^2}{2\mu_0}. \quad (2.547)$$

This suggests that the energy density of a general electromagnetic field is

$$U = \frac{\epsilon_0 E^2}{2} + \frac{B^2}{2\mu_0}. \quad (2.548)$$

Let us now demonstrate that Maxwell's equations conserve energy. We have already come across one conservation law in electromagnetism; namely,

$$\frac{\partial \rho}{\partial t} + \nabla \cdot \mathbf{j} = 0. \quad (2.549)$$

(See Equation (2.473).] The previous expression described the conservation of electric charge. Thus, integrating over some volume V , bounded by a surface S , and making use of the divergence theorem (see Section A.20), we obtain

$$-\frac{\partial}{\partial t} \int_V \rho dV = \oint_S \mathbf{j} \cdot d\mathbf{S}. \quad (2.550)$$

In other words, the rate of decrease of the electric charge contained in volume V is equal to the net flux of charge out of surface S . This suggests that an energy conservation law for electromagnetic fields should have the form

$$-\frac{\partial}{\partial t} \int_V U dV = \oint_S \mathbf{u} \cdot d\mathbf{S}. \quad (2.551)$$

Here, U is the energy density of the electromagnetic field, and \mathbf{u} is the flux of electromagnetic energy (i.e., energy $|\mathbf{u}|$ per unit time, per unit cross-sectional area, passes a given point in the direction of \mathbf{u}). According to the previous equation, the rate of decrease of the electromagnetic energy in volume V is equal to the net flux of electromagnetic energy out of surface S .

Equation (2.551) is actually incomplete, because electromagnetic fields can gain or lose energy by interacting with matter. We need to incorporate this fact into our analysis. We saw earlier (see Section 2.1.11) that the rate of heat dissipation per unit volume in a conductor (the so-called ohmic heating rate) is $\mathbf{E} \cdot \mathbf{j}$. This energy is extracted from electromagnetic fields, so the rate of energy loss of the fields in a volume V due to interaction with matter is $\int_V \mathbf{E} \cdot \mathbf{j} dV$. Thus, Equation (2.551) generalizes to give

$$-\frac{\partial}{\partial t} \int_V U dV = \oint_S \mathbf{u} \cdot d\mathbf{S} + \int_V \mathbf{E} \cdot \mathbf{j} dV. \quad (2.552)$$

From the divergence theorem (see Section A.20), the previous equation is equivalent to

$$\frac{\partial U}{\partial t} + \nabla \cdot \mathbf{u} = -\mathbf{E} \cdot \mathbf{j}. \quad (2.553)$$

Let us now see if we can derive an expression of this form from Maxwell's equations.

We start from the differential form of Ampère's law (including the displacement current), (2.487):

$$\nabla \times \mathbf{B} = \mu_0 \mathbf{j} + \epsilon_0 \mu_0 \frac{\partial \mathbf{E}}{\partial t}. \quad (2.554)$$

The scalar product of the electric field with this equation yields

$$-\mathbf{E} \cdot \mathbf{j} = -\frac{\mathbf{E} \cdot \nabla \times \mathbf{B}}{\mu_0} + \epsilon_0 \mathbf{E} \cdot \frac{\partial \mathbf{E}}{\partial t}. \quad (2.555)$$

The previous expression can be rewritten

$$-\mathbf{E} \cdot \mathbf{j} = -\frac{\mathbf{E} \cdot \nabla \times \mathbf{B}}{\mu_0} + \frac{\partial}{\partial t} \left(\frac{\epsilon_0 E^2}{2} \right). \quad (2.556)$$

However (see Section A.24),

$$\nabla \cdot (\mathbf{E} \times \mathbf{B}) \equiv \mathbf{B} \cdot \nabla \times \mathbf{E} - \mathbf{E} \cdot \nabla \times \mathbf{B}, \quad (2.557)$$

so

$$-\mathbf{E} \cdot \mathbf{j} = \nabla \cdot \left(\frac{\mathbf{E} \times \mathbf{B}}{\mu_0} \right) - \frac{\mathbf{B} \cdot \nabla \times \mathbf{E}}{\mu_0} + \frac{\partial}{\partial t} \left(\frac{\epsilon_0 E^2}{2} \right). \quad (2.558)$$

The differential form of Faraday's law, (2.486), yields

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}, \quad (2.559)$$

so

$$-\mathbf{E} \cdot \mathbf{j} = \nabla \cdot \left(\frac{\mathbf{E} \times \mathbf{B}}{\mu_0} \right) + \mu_0^{-1} \mathbf{B} \cdot \frac{\partial \mathbf{B}}{\partial t} + \frac{\partial}{\partial t} \left(\frac{\epsilon_0 E^2}{2} \right), \quad (2.560)$$

which can be rewritten as

$$-\mathbf{E} \cdot \mathbf{j} = \nabla \cdot \left(\frac{\mathbf{E} \times \mathbf{B}}{\mu_0} \right) + \frac{\partial}{\partial t} \left(\frac{\epsilon_0 E^2}{2} + \frac{B^2}{2\mu_0} \right). \quad (2.561)$$

Thus, we obtain the desired conservation law,

$$\frac{\partial U}{\partial t} + \nabla \cdot \mathbf{u} = -\mathbf{E} \cdot \mathbf{j}, \quad (2.562)$$

where

$$U = \frac{\epsilon_0 E^2}{2} + \frac{B^2}{2\mu_0} \quad (2.563)$$

is the electromagnetic energy density, and

$$\mathbf{u} = \frac{\mathbf{E} \times \mathbf{B}}{\mu_0} \quad (2.564)$$

is the electromagnetic energy flux. The latter quantity is usually called the *Poynting flux*, after John Poynting who derived it in 1884.

Let us see whether our expression for the electromagnetic energy flux makes physical sense. We know that if we stand in the Sun then we get hot. This occurs because we absorb electromagnetic radiation emitted by the Sun. So, radiation must transport energy. The electric and magnetic fields in electromagnetic radiation are mutually perpendicular, and are also perpendicular to the direction of propagation, $\hat{\mathbf{k}}$ (which is a unit vector). Furthermore, $B = E/c$. (See Section 2.4.4.)

Equation (2.521) can easily be transformed into the following relation between the electric and magnetic fields of an electromagnetic wave:

$$\mathbf{E} \times \mathbf{B} = \frac{E^2}{c} \hat{\mathbf{k}}. \quad (2.565)$$

Thus, the Poynting flux for electromagnetic radiation is

$$\mathbf{u} = \frac{E^2}{\mu_0 c} \hat{\mathbf{k}} = \epsilon_0 c E^2 \hat{\mathbf{k}}. \quad (2.566)$$

The previous expression states that electromagnetic waves transport energy along their direction of propagation, which seems to make sense.

The energy density of electromagnetic radiation is

$$U = \frac{\epsilon_0 E^2}{2} + \frac{B^2}{2\mu_0} = \frac{\epsilon_0 E^2}{2} + \frac{E^2}{2\mu_0 c^2} = \epsilon_0 E^2, \quad (2.567)$$

where use has been made of $B = E/c$, and $c = 1/\sqrt{\epsilon_0 \mu_0}$. Note that the electric and magnetic components of an electromagnetic wave have equal energy densities. Because electromagnetic waves travel at the speed of light in vacuum, we would expect the energy flux through one square meter in one second to equal the energy contained in a volume of length c , and unit cross-sectional area; that is, c multiplied by the electromagnetic energy density. Thus,

$$|\mathbf{u}| = c U = \epsilon_0 c E^2, \quad (2.568)$$

which is in accordance with the previous two equations.

2.4.6 Electromagnetic Momentum

We have seen that electromagnetic waves carry energy. It turns out that they also carry momentum. Consider the following argument, due to Einstein. Suppose that we have a railroad car of mass M and length L that is free to move in one dimension. See Figure 2.42. Suppose that electromagnetic radiation of total energy E is emitted from one end of the car, propagates along the length of the car, and is then absorbed at the other end. The effective mass of this radiation is $m = E/c^2$ (from Einstein's famous relation $E = m c^2$). (See Section 3.3.4.) At first sight, the process described previously appears to cause the center of mass of the system to spontaneously shift. This violates the law of momentum conservation (assuming the railway car is subject to no net horizontal external force). (See Section 1.4.4.) The only way in which the center of mass of the system can remain stationary is if the railway car moves in the opposite direction to the direction of propagation of the radiation. In fact, if the car moves by a distance x then the center of mass of the system is the same before and after the radiation pulse provided that

$$M x = m L = \frac{E}{c^2} L. \quad (2.569)$$

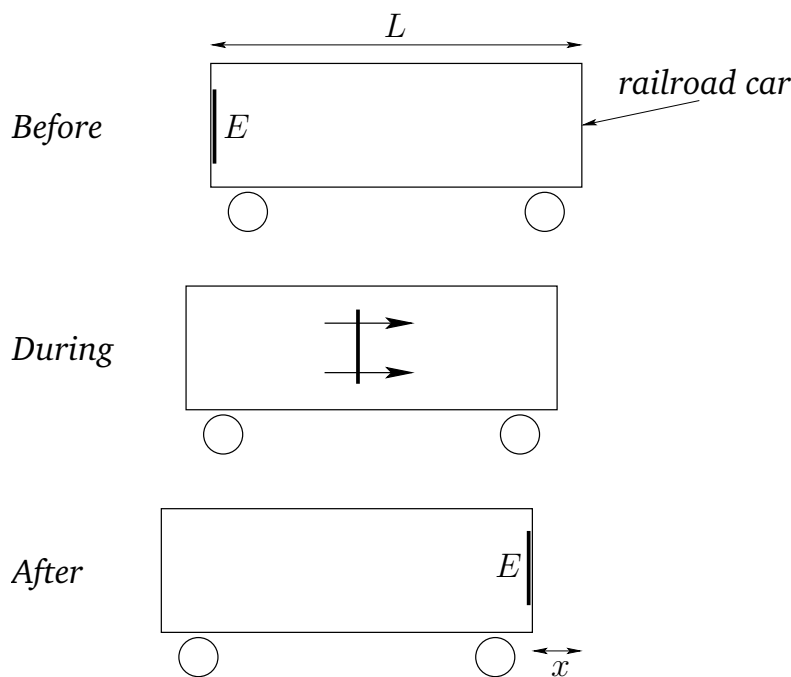


Figure 2.42: Einstein's thought experiment regarding electromagnetic momentum.

Incidentally, it is assumed that $m \ll M$ in this derivation.

But, what actually causes the car to move? If the radiation possesses momentum, p , then the car will recoil with the same momentum when the radiation is emitted. When the radiation hits the other end of the car then the car acquires momentum p in the opposite direction, which halts the motion. The time of flight of the radiation is L/c . So, the distance traveled by a mass M with momentum p in this time is

$$x = vt = \frac{p L}{M c}, \quad (2.570)$$

giving

$$p = M x \frac{c}{L} = \frac{E}{c}. \quad (2.571)$$

Thus, the momentum carried by electromagnetic radiation is equal to its energy divided by the speed of light. The same result can be obtained from the well-known relativistic formula

$$E^2 = p^2 c^2 + m^2 c^4 \quad (2.572)$$

relating the energy E , momentum p , and mass m of a particle. (See Section 3.3.5.) According to quantum theory, electromagnetic radiation is made up of massless particles called *photons*. (See Section 3.3.8.) Thus,

$$p = \frac{E}{c} \quad (2.573)$$

for individual photons, so the same must be true of electromagnetic radiation as a whole.

It follows from Equation (2.571) that the momentum density, g , of electromagnetic radiation is equal to its energy density divided by c , so that

$$g = \frac{U}{c} = \frac{\epsilon_0 E^2}{c} = \frac{|\mathbf{u}|}{c^2}, \quad (2.574)$$

where use has been made of Equations (2.566) and (2.567). It is reasonable to suppose that the momentum is directed along the direction of the energy flow (this is obviously the case for photons), so the vector momentum density (which gives the direction, as well as the magnitude, of the momentum per unit volume) of electromagnetic radiation is

$$\mathbf{g} = \frac{\mathbf{u}}{c^2}. \quad (2.575)$$

Thus, the momentum density of electromagnetic fields is equal to the associated energy flux divided by c^2 .

Of course, the electric field associated with an electromagnetic wave oscillates rapidly in time, which implies that the previous expressions for the energy density, energy flux, and momentum density of electromagnetic radiation also oscillate rapidly. It is convenient to average over many periods of the oscillation (this average is denoted $\langle \rangle$). Thus, from Equations (2.566), (2.567), and (2.575),

$$\langle U \rangle = \frac{\epsilon_0 E_0^2}{2}, \quad (2.576)$$

$$\langle \mathbf{u} \rangle = \frac{c \epsilon_0 E_0^2}{2} \hat{\mathbf{k}} = c \langle U \rangle \hat{\mathbf{k}}, \quad (2.577)$$

$$\langle \mathbf{g} \rangle = \frac{\epsilon_0 E_0^2}{2c} \hat{\mathbf{k}} = \frac{\langle U \rangle}{c} \hat{\mathbf{k}}, \quad (2.578)$$

where the factor $1/2$ comes from averaging $\cos^2(\omega t)$. Here, E_0 is the peak amplitude of the electric field associated with the wave.

If electromagnetic radiation possesses momentum then it must exert a force on bodies that absorb (or emit) radiation. Suppose that a body is placed in a beam of perfectly collimated radiation, that it completely absorbs. The amount of momentum absorbed per unit time, per unit cross-sectional area, is simply the amount of momentum contained in a volume of length c , and unit cross-sectional area; that is, c multiplied by the momentum density, g . An absorbed momentum per unit time, per unit area, is equivalent to a pressure. In other words, the radiation exerts a pressure $c g$ on the body. Thus, the *radiation pressure* is given by

$$p = \frac{\epsilon_0 E^2}{2} = \langle U \rangle. \quad (2.579)$$

So, the pressure exerted by collimated electromagnetic radiation is equal to its average energy density.

Consider a cavity filled with electromagnetic radiation. What is the radiation pressure exerted on the walls? In this situation, the radiation propagates in all directions with equal probability.

Consider radiation propagating at an angle θ to the local normal to the wall. The amount of such radiation hitting the wall per unit time, per unit area, is proportional to $\cos \theta$. Moreover, the component of momentum normal to the wall that the radiation carries is also proportional to $\cos \theta$. Thus, the pressure exerted on the wall is the same as in Equation (2.579), except that it is weighted by the average of $\cos^2 \theta$ over all solid angles, in order to take into account the fact that obliquely propagating radiation exerts a pressure that is $\cos^2 \theta$ times that of normal radiation. The average of $\cos^2 \theta$ over all solid angles is $1/3$, so for isotropic radiation

$$p = \frac{\langle U \rangle}{3}. \quad (2.580)$$

Clearly, the pressure exerted by isotropic radiation is one third of its average energy density. Radiation pressure is unimportant in the Sun, but makes a significant contribution to counteracting gravitational collapse in larger, hotter stars.

The power incident on the surface of the Earth, due to radiation emitted by the Sun, is about 1300 W m^{-2} . So, what is the radiation pressure? Because,

$$\langle |\mathbf{u}| \rangle = c \langle U \rangle = 1300 \text{ W m}^{-2}, \quad (2.581)$$

then

$$p = \langle U \rangle \simeq 4 \times 10^{-6} \text{ N m}^{-2}. \quad (2.582)$$

Here, the radiation is assumed to be perfectly collimated. Thus, the radiation pressure exerted on the Earth is minuscule (for comparison, the pressure of the atmosphere is about 10^5 N m^{-2}). Nevertheless, this small pressure due to radiation is important in outer space, because it is responsible for continuously sweeping dust particles out of the solar system. It is quite common for comets to exhibit two separate tails. One, known as the *gas tail*, consists of ionized gas, and is swept along by the solar wind (a stream of charged particles and magnetic field-lines emitted by the Sun). The other, known as the *dust tail*, consists of uncharged dust particles, and is swept radially outward (because light travels in straight-lines) from the Sun by radiation pressure. Two separate tails are observed if the local direction of the solar wind is not radially outward from the Sun (which is quite often the case).

The radiation pressure from sunlight is very weak. However, that produced by laser beams can be enormous (far higher than any conventional pressure which has ever been produced in a laboratory). For instance, the lasers used in inertial confinement fusion (e.g., the National Ignition Facility at Lawrence Livermore National Laboratory) typically have energy fluxes of 10^{18} W m^{-2} . This translates to a radiation pressure of about 10^4 atmospheres.

Chapter 3

Special Relativity

3.1 Experimental Basis of Special Relativity

3.1.1 Sound Waves in a Gas

A *sound wave* in a gas is a longitudinal disturbance of the gas's pressure and density that propagates at the fixed speed

$$c = \sqrt{\frac{\gamma p}{\rho}}, \quad (3.1)$$

(See Section 5.2.9.) Here, γ is the gas's ratio of specific heats (which is approximately 1.4 for the atmosphere), p the gas's undisturbed pressure, and ρ the gas's undisturbed mass density. Note that a sound wave is a *non-dispersive* wave, which means that a transient wave pulse propagates at the same speed as an infinite wave train. (See Section 4.2.6.) However, a sound wave only propagates at the speed (3.1) in the rest frame of the gas.

Let \mathbf{c} be the phase velocity of a sound wave in a stationary frame of reference in which the gas is at rest. Thus, $|\mathbf{c}| = c$, is the speed of sound, (3.1). Consider a moving frame of reference that moves at constant velocity \mathbf{v} , where $v < c$, with respect to the stationary frame. Incidentally, if the stationary frame is inertial then so is the moving frame. (See Section 1.5.4.) The gas appears to flow with uniform velocity $-\mathbf{v}$ in the moving reference frame. Furthermore, it is an experimentally verified fact that the sound wave appears to propagate with the phase velocity

$$\mathbf{c}' = \mathbf{c} - \mathbf{v} \quad (3.2)$$

in the moving frame. Note that, in general, both the speed and the direction of the sound wave are different in the stationary and the moving frames. Note, further, that the previous equation is a direct consequence of the Galilean transformation, (1.106)–(1.108). (See Section 3.2.6.) The previous equation yields

$$c' = (c^2 - 2\mathbf{c}' \cdot \mathbf{v} - v^2)^{1/2}. \quad (3.3)$$

Hence, we deduce, from the previous two equations, that if the sound wave propagates in the same direction as \mathbf{v} in the moving frame then it propagates at the speed

$$u_- = c - v, \quad (3.4)$$

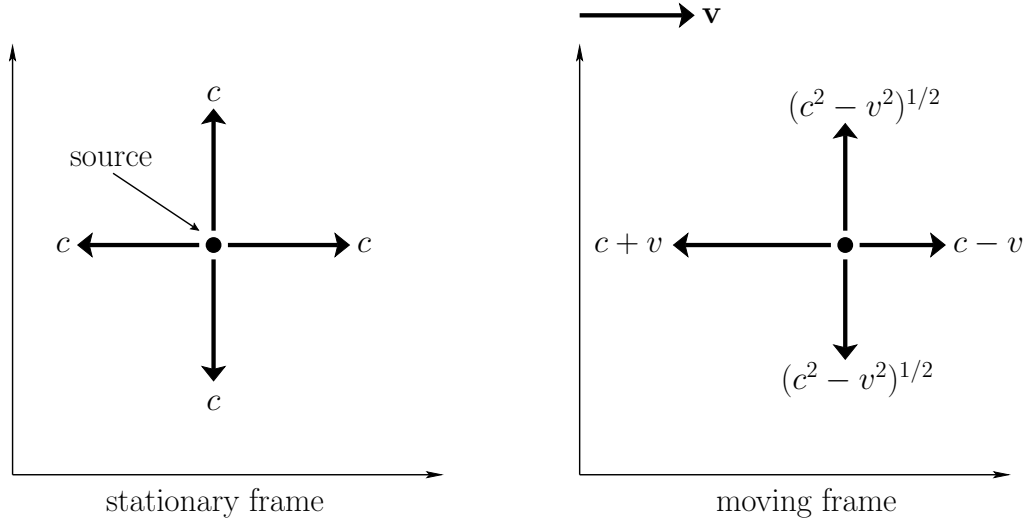


Figure 3.1: Sound waves propagating in a stationary and in a moving reference frame.

but if the sound wave propagates in the opposite direction to \mathbf{v} in the moving frame then it propagates at the speed

$$u_+ = c + v \quad (3.5)$$

and, finally, if the sound wave propagates in a direction perpendicular to \mathbf{v} in the moving frame then it propagates at the speed

$$u_{\perp} = (c^2 - v^2)^{1/2}. \quad (3.6)$$

Of course, the sound wave propagates at the speed c in all directions in the stationary frame. These ideas are illustrated in Figure 3.1.

We could imagine performing an experiment in the moving reference frame in order to measure its velocity, \mathbf{v} , with respect to the stationary frame. See Figure 3.2. Suppose that we have a sound wave source that emits a transient sound wave pulse isotropically in all directions. Suppose that we place two small sound-wave reflectors (which are stationary in the moving frame) at equal distances l_0 from the source. The displacement of the first reflector from the source is in the direction of \mathbf{v} , whereas the displacement of the second reflector is in a direction that is perpendicular to \mathbf{v} . The travel time of the pulse from the source to the first reflector, and back again, is

$$t_1 = \frac{l_0}{u_-} + \frac{l_0}{u_+} = \frac{l_0}{c - v} + \frac{l_0}{c + v} = \frac{2 l_0 c}{c^2 - v^2} \approx \frac{2 l_0}{c} \left(1 + \frac{v^2}{c^2} \right), \quad (3.7)$$

where we have assumed that $v \ll c$. The travel time of the pulse from the source to the second reflector, and back again, is

$$t_2 = \frac{2 l_0}{u_{\perp}} = \frac{2 l_0}{(c^2 - v^2)^{1/2}} \approx \frac{2 l_0}{c} \left(1 + \frac{v^2}{2 c^2} \right). \quad (3.8)$$

Thus, if we measure the two travel times, and take the difference between them, then we obtain

$$t_1 - t_2 = \frac{l_0 v^2}{c^3}. \quad (3.9)$$

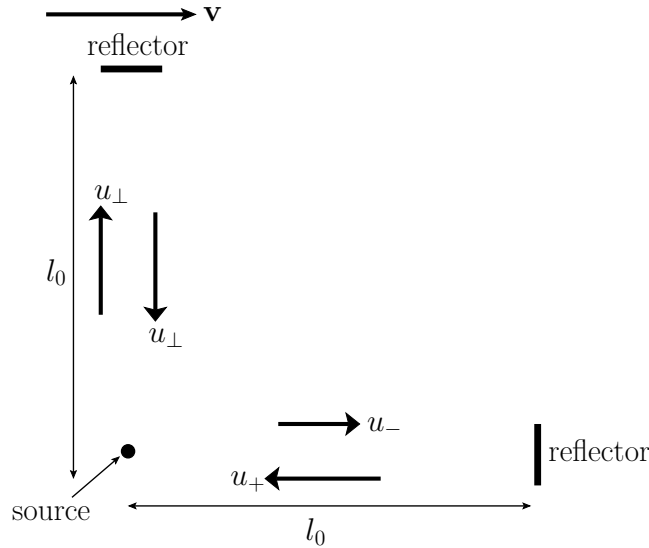


Figure 3.2: Experiment to detect motion of moving inertial reference frame.

Hence, assuming that we know l_0 and c , we can determine v . We can also determine the direction of v because the time difference is maximized when the two legs of the apparatus shown in Figure 3.2 are aligned parallel and perpendicular to this direction.

3.1.2 Light Waves in a Vacuum

A *light wave* is a transverse disturbance of electric and magnetic fields that is able to propagate through a vacuum (unlike a sound wave), and does so at the fixed speed

$$c = \frac{1}{\sqrt{\epsilon_0 \mu_0}}, \quad (3.10)$$

where $\epsilon_0 = 8.854 \times 10^{-12} \text{ F m}^{-1}$ is the *electrical permittivity of free space*, and $\mu_0 = 4\pi \times 10^{-7} \text{ H m}^{-1}$ the *magnetic permeability of free space*. (See Section 2.4.4.) It follows that

$$c = \frac{1}{[(8.854 \times 10^{-12})(4\pi \times 10^{-7})]^{1/2}} = 2.998 \times 10^8 \text{ m s}^{-1}. \quad (3.11)$$

Note that ϵ_0 and μ_0 can be determined from simple experiments involving measurements of the forces exerted by electric charges and current loops on one another.

The classical theory of electromagnetism (i.e., Maxwell's equations) does not explicitly mention a medium through which electromagnetic disturbances propagate. (See Section 2.4.2.) Nevertheless, prior to the 20th century, most physicists assumed that such a medium existed, because they could not conceive of a wave that propagated in the absence of a medium. The medium in question was known as the *aether* (from the ancient Greek *αἰθήρ*, which is the fifth element of Aristotelian philosophy), and was thought to permeate all space, including vacuums. Thus, by

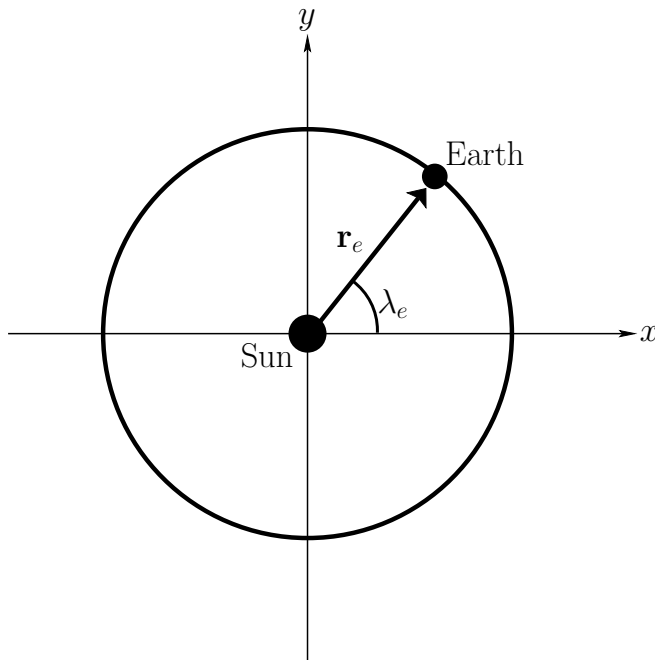


Figure 3.3: Orbital motion of the Earth about the Sun.

analogy with a sound wave, the phase velocity of a light wave in a frame of reference moving at fixed velocity \mathbf{v} with respect to the rest frame of the aether was assumed to be

$$\mathbf{c}' = \mathbf{c} - \mathbf{v}, \quad (3.12)$$

where \mathbf{c} is the phase velocity of the light wave in the rest frame of the aether, and $|\mathbf{c}| = c$ is the speed of light, (3.10).

3.1.3 Aberration of Starlight and Stellar Parallax

Prior to the 20th century, one of the strongest arguments in favor of the existence of the aether was thought to be the *aberration of light*. This is a phenomenon that produces an apparent motion of distant stars about their true positions, due to a combination of the finite velocity of light and the Earth's orbital motion about the Sun. Aberration is closely related to another phenomenon, known as *parallax*, that also produces an apparent motion of distant stars about their true positions; in this case, due to the Earth's shifting position about the Sun. It is convenient to discuss these two effects together.

The Earth moves around the Sun in a planar orbit whose plane includes the Sun. (See Section 1.9.2.) The orbit is approximately circular in shape, and has a radius $a_e = 1.496 \times 10^{11}$ m. (See Sections 1.9.6, and Table 1.4.) The plane that contains the Earth's orbit is known as the *ecliptic plane*. (See Section 1.10.6.) Let us set up a Cartesian coordinate system in the ecliptic plane whose origin coincides with the Sun, whose z -axis is directed toward the northern ecliptic pole (i.e., the direction that is normal to the ecliptic plane in a northern sense), and whose x -axis is directed

toward the *vernal equinox* (i.e., the point in the sky at which the Sun annually passes through the projection of the Earth's equatorial plane in a northward sense). See Figure 3.3. The angle, λ_e , shown in the figure, is known as the Earth's *ecliptic longitude*, and serves to locate the Earth on its orbit. Let \mathbf{r}_e be the displacement of the Earth from the Sun. It is clear from simple geometry that the components of \mathbf{r}_e are

$$\mathbf{r}_e = a_e (\cos \lambda_e, \sin \lambda_e, 0). \quad (3.13)$$

Thus, the Earth's orbital velocity becomes

$$\mathbf{v}_e = a_e \frac{d\lambda_e}{dt} (-\sin \lambda_e, \cos \lambda_e, 0), \quad (3.14)$$

where

$$\frac{d\lambda_e}{dt} = \left(\frac{G M_s}{a_e^3} \right)^{1/2} \quad (3.15)$$

is the Earth's mean orbital angular velocity about the Sun. (See Section 1.9.7.) Here, $M_s = 1.989 \times 10^{30}$ kg is the Sun's mass. Let λ_s be the apparent ecliptic longitude of the Sun, as seen on the Earth. It is clear that $\lambda_s = \pi - \lambda_e$. Hence, we deduce that

$$\mathbf{r}_e = -a_e \mathbf{e}_r, \quad (3.16)$$

$$\mathbf{v}_e = -v_e \mathbf{e}_\theta, \quad (3.17)$$

$$\mathbf{e}_r = (\cos \lambda_s, \sin \lambda_s, 0), \quad (3.18)$$

$$\mathbf{e}_\theta = (-\sin \lambda_s, \cos \lambda_s, 0). \quad (3.19)$$

Here, \mathbf{e}_r is a unit vector directed from the Earth to the Sun, whereas \mathbf{e}_θ is a unit vector that is parallel to the Sun's apparent orbital velocity about the Earth. Finally,

$$v_e = \left(\frac{G M_s}{a_e} \right)^{1/2} = 2.977 \times 10^4 \text{ m s}^{-1} \quad (3.20)$$

is the Earth's mean orbital velocity about the Sun.

Suppose that a light ray from a distant star is observed on the Earth. Let the phase velocity of the ray in the aether rest frame, in which the Sun is assumed to be stationary, be \mathbf{c} , where $|\mathbf{c}| = c$ is the speed of light in vacuum. Because the Earth is actually moving with respect to the aether rest frame, the observed phase velocity of the light ray is

$$\mathbf{c}' = \mathbf{c} - \mathbf{v}_e. \quad (3.21)$$

[See Equation (3.12).] Let θ be the angle subtended between \mathbf{v}_e and $-\mathbf{c}$, and let θ' be the angle subtended between \mathbf{v}_e and $-\mathbf{c}'$. See Figure 3.4. Thus, θ corresponds to the true angular location of the star (i.e., the location seen by an observer in the Sun's rest frame), whereas θ' corresponds to the apparent location of the star seen on the moving Earth. Simple trigonometry reveals that

$$\frac{\sin(\theta - \theta')}{v_e} = \frac{\sin \theta'}{c}. \quad (3.22)$$

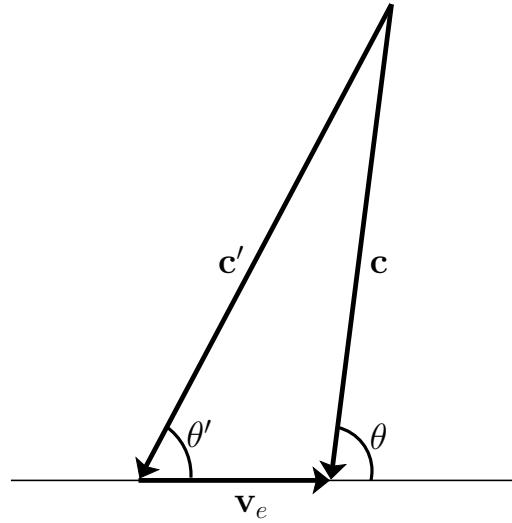


Figure 3.4: Aberration of starlight.

However, $\sin(\theta' - \theta) = \sin \theta' \cos \theta - \cos \theta' \sin \theta$, so we get

$$\tan \theta' = \frac{\sin \theta}{\cos \theta + \kappa}, \quad (3.23)$$

where

$$\kappa = \frac{v_e}{c} = 9.930 \times 10^{-5} \quad (3.24)$$

is known as the *constant of aberration*.

Let us write

$$\mathbf{c} = -c \mathbf{n}, \quad (3.25)$$

$$\mathbf{c}' = -c' \mathbf{n}', \quad (3.26)$$

where the unit vectors \mathbf{n} and \mathbf{n}' are directed toward the true position of the star (in the Sun's rest frame), and the apparent position seen on the moving Earth, respectively. It is clear from Equations (3.17), (3.21), and (3.24) that

$$\mathbf{n}' = \frac{\mathbf{n} - \kappa \mathbf{e}_\theta}{|\mathbf{n} - \kappa \mathbf{e}_\theta|} \simeq \mathbf{n} - \kappa [\mathbf{e}_\theta - (\mathbf{n} \cdot \mathbf{e}_\theta) \mathbf{n}], \quad (3.27)$$

to first order in κ . Let us write

$$\mathbf{n} = (\cos \beta \cos \lambda, \cos \beta \sin \lambda, \sin \beta), \quad (3.28)$$

$$\mathbf{n}' = (\cos \beta' \cos \lambda', \cos \beta' \sin \lambda', \sin \beta'). \quad (3.29)$$

Here, β and λ are the true *ecliptic latitude* and *ecliptic longitude* of the star, respectively, whereas β' and λ' are the apparent latitude and longitude seen on the moving Earth. (Ecliptic latitude and

longitude parameterize position on the celestial sphere, and are similar to terrestrial latitude and longitude, except that the equator corresponds to the Earth's orbital plane, and ecliptic longitude increases in the opposite direction to terrestrial longitude.) It is clear from Equations (3.19) and (3.28) that

$$\mathbf{n} \cdot \mathbf{e}_\theta = -\cos\beta \sin(\lambda_s - \lambda). \quad (3.30)$$

Hence, Equations (3.27)–(3.29) yield

$$\cos\beta' \cos\lambda' = \cos\beta \cos\lambda - \kappa \cos^2\beta \cos\lambda \sin(\lambda_s - \lambda) + \kappa \sin\lambda_s, \quad (3.31)$$

$$\cos\beta' \sin\lambda' = \cos\beta \sin\lambda - \kappa \cos^2\beta \sin\lambda \sin(\lambda_s - \lambda) - \kappa \cos\lambda_s, \quad (3.32)$$

$$\sin\beta' = \sin\beta - \kappa \cos\beta \sin\beta \sin(\lambda_s - \lambda). \quad (3.33)$$

Equations (3.31) and (3.32) can be combined to give

$$\cos\beta' \sin(\lambda' - \lambda) = -\kappa \cos(\lambda_s - \lambda). \quad (3.34)$$

Finally, writing $\beta' = \beta + \delta\beta$ and $\lambda' = \lambda + \delta\lambda$, Equations (3.33) and (3.34) yield

$$\delta\beta = -\kappa \sin\beta \sin(\lambda_s - \lambda), \quad (3.35)$$

$$\delta\lambda = -\frac{\kappa}{\cos\beta} \cos(\lambda_s - \lambda), \quad (3.36)$$

to first order in κ . If $x = \delta\lambda \cos\beta$ represents angular displacement on the celestial sphere in a direction parallel to the ecliptic plane (in the sense of the Sun's apparent motion with respect to the stars), and $y = \delta\beta$ represents angular displacement in a direction perpendicular to the ecliptic (in a northern sense), then the previous two equations give

$$x = -\kappa \cos(\lambda_s - \lambda), \quad (3.37)$$

$$y = -\kappa \sin\beta \sin(\lambda_s - \lambda). \quad (3.38)$$

This is clearly the parametric equation of an ellipse. Hence, we deduce that, as a consequence of the aberration of light, during the course of a year, our star appears to describe an ellipse on the celestial sphere. The major radius, κ , is parallel to the ecliptic plane, whereas the minor radius, $\kappa \sin\beta$, is perpendicular to the ecliptic. The angular displacement of the star from its mean position is greatest when $\lambda_s - \lambda = 0^\circ$, or 180° (i.e., when the ecliptic longitude of the star matches that of the Sun, or differs from it by 180° , which maximizes the Earth's transverse velocity with respect to the star). The magnitude of the greatest angular displacement, κ , is 20.48 arc seconds. This is about the same as the angular size of Saturn's disk.

Let us now consider parallax. Let \mathbf{d} be the displacement of the Sun from a distant star, let \mathbf{d}' be the corresponding displacement of the Earth, and let \mathbf{r}_e be the displacement of the Earth from the Sun. It is evident that

$$\mathbf{d}' = \mathbf{d} + \mathbf{r}_e. \quad (3.39)$$

See Figure 3.5. Let θ be the angle subtended between \mathbf{r}_e and $-\mathbf{d}$, and let θ' be the angle subtended between \mathbf{r}_e and $-\mathbf{d}'$. Thus, θ corresponds to the true location of the star (i.e., the location seen by

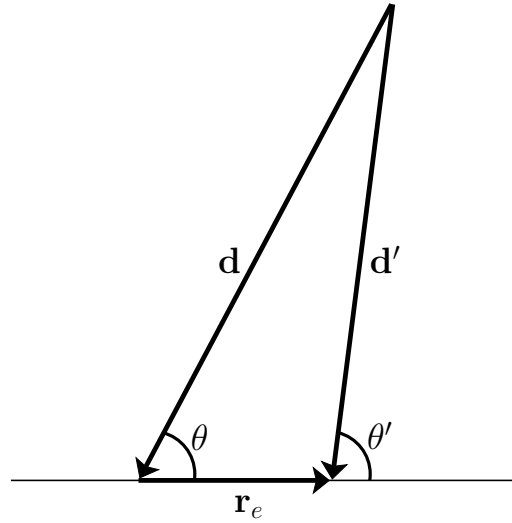


Figure 3.5: Stellar parallax.

an observer on the Sun), whereas θ' corresponds to the apparent location of the star seen on the displaced Earth. Simple trigonometry reveals that

$$\frac{\sin(\theta' - \theta)}{a_e} = \frac{\sin \theta'}{d}. \quad (3.40)$$

However, $\sin(\theta' - \theta) = \sin \theta' \cos \theta - \cos \theta' \sin \theta$, so we get

$$\tan \theta' = \frac{\sin \theta}{\cos \theta - \Pi}, \quad (3.41)$$

where

$$\Pi = \frac{a_e}{d} \quad (3.42)$$

is known as the star's *parallax*. If we measure the star's distance from the Sun, d , in units of *parsecs* (pc) then, by definition, the parallax in arc seconds is

$$\Pi = \frac{1}{d}. \quad (3.43)$$

It follows that

$$1 \text{ pc} = \frac{60 \times 60 \times 180 a_e}{\pi} = 3.087 \times 10^{16} \text{ m}. \quad (3.44)$$

Given that the nearest star to the Sun, Proxima Centauri, is 1.302 parsecs away, it is clear that all stellar parallaxes are less than 1 arc second. This implies that stellar aberration is a much larger effect than stellar parallax.

We can write $\mathbf{d} = -d \mathbf{n}$ and $\mathbf{d}' = -d' \mathbf{n}'$. Making use of very similar analysis to that used to calculate aberration, we obtain

$$\mathbf{n}' = \frac{\mathbf{n} + \Pi \mathbf{e}_r}{|\mathbf{n} + \Pi \mathbf{e}_r|} \approx \mathbf{n} + \Pi [\mathbf{e}_r - (\mathbf{n} \cdot \mathbf{e}_r) \mathbf{n}], \quad (3.45)$$

$$\mathbf{n} \cdot \mathbf{e}_r = \cos\beta \cos(\lambda_s - \lambda), \quad (3.46)$$

$$\cos\beta' \cos\lambda' = \cos\beta \cos\lambda - \Pi \cos^2\beta \cos\lambda \cos(\lambda_s - \lambda) + \Pi \cos\lambda_s, \quad (3.47)$$

$$\cos\beta' \sin\lambda' = \cos\beta \sin\lambda - \Pi \cos^2\beta \sin\lambda \cos(\lambda_s - \lambda) + \Pi \sin\lambda_s, \quad (3.48)$$

$$\sin\beta' = \sin\beta - \Pi \cos\beta \sin\beta \sin(\lambda_s - \lambda). \quad (3.49)$$

Equations (3.47) and (3.48) can be combined to give

$$\cos\beta' \sin(\lambda' - \lambda) = \Pi \sin(\lambda_s - \lambda). \quad (3.50)$$

Hence, Equations (3.49) and (3.50) yield

$$\delta\beta = -\Pi \sin\beta \cos(\lambda_s - \lambda), \quad (3.51)$$

$$\delta\lambda = \frac{\Pi}{\cos\beta} \sin(\lambda_s - \lambda), \quad (3.52)$$

to first order in Π . If we again let $x = \delta\lambda \cos\beta$ represent angular displacement on the celestial sphere in a direction parallel to the ecliptic plane, and $y = \delta\beta$ represent angular displacement in a direction perpendicular to the ecliptic, then the previous two equations give

$$x = \Pi \sin(\lambda_s - \lambda), \quad (3.53)$$

$$y = -\Pi \sin\beta \cos(\lambda_s - \lambda). \quad (3.54)$$

This is again the parametric equation of an ellipse. Hence, we deduce that, as a consequence parallax, during the course of a year, our star appears to describe an ellipse on the celestial sphere. The major radius, Π , is parallel to the ecliptic plane, whereas the minor radius, $\Pi \sin\beta$, is perpendicular to the ecliptic. The angular displacement of the star from its mean position is greatest when $\lambda_s - \lambda = 90^\circ$, or 270° (i.e., when the ecliptic longitude of the star differs from that of the Sun by 90° or 270° , which maximizes the Earth's transverse displacement with respect to the star). The greatest angular displacement, Π , when measured in arc seconds, is equal to one over the distance of the star from the Sun measured in parsecs. [See Equation (3.43).]

Between 1725 and 1727, the astronomers James Bradley and Samuel Molyneux measured the position of the circumpolar star γ Draconis in an attempt to observe its parallax. They found that the angular position of the star underwent small annual variations, but that the deviation from the mean was greatest when the ecliptic longitude of the Sun matched that of the star, which is not the behavior expected from parallax. In 1728, Bradley correctly explained the observed variations in terms of the aberration of light. Note that this explanation depends crucially on the fact that the speed of light observed in a frame of reference that moves with respect to the rest frame of the aether is different to that observed in the rest frame. Incidentally, stellar parallax is so small an effect that it was not successfully measured until 1838, when Frederick Bessel measured the parallax of the star 61 Cygni.

3.1.4 Fizeau and Airy Experiments

Light can also propagate through transparent dielectric media, such as air and water, but does so at the reduced phase velocity, c/n , where c is the speed of light in vacuum, and n is the medium's refractive index. Prior to the mid-nineteenth century, it was suppose that the aether was entrained by a moving medium, so that if v is the speed of the medium then the phase velocity of light is

$$u_+ = \frac{c}{n} + v \quad (3.55)$$

when it propagates in the same direction as the medium, and

$$u_- = \frac{c}{n} - v \quad (3.56)$$

when it propagates in the opposite direction [c.f. Equations (3.4) and (3.5).] However, in 1951, Hippolyte Fizeau measured the speed of light in moving water, and found that

$$u_+ = \frac{c}{n} + v \left(1 - \frac{1}{n^2} \right), \quad (3.57)$$

$$u_- = \frac{c}{n} - v \left(1 - \frac{1}{n^2} \right). \quad (3.58)$$

He concluded that the aether is only partially dragged by a moving medium. In particular, air, which has a refractive index of 1.0003, hardly drags the aether at all. This is a good thing because a strong aether drag through air would contradict Bradley's explanation of stellar aberration. On the other hand, water, which has a refractive index of 1.33, drags the aether by 43% of its velocity. Unfortunately, in 1871, George Airy demonstrated that measured stellar aberration is the same when the observing telescope is filled with water as when it is filled with air. This completely contradicts the partial aether drag hypothesis.

3.1.5 Michelson-Morley Experiment

The aberration of stellar light can be thought of as an indirect measurement of the velocity of the Earth with respect to the aether. In 1887, Albert Michelson and Edward Morley attempted to measure this velocity directly in the laboratory. Their apparatus was an optical version of the hypothetical sound-wave experiment shown in Figure 3.2. The apparatus sends light through a half-silvered mirror that is used to split it into two beams that travel at right angles to one another. The two beams travel out to the ends of long arms of equal length where they are reflected back into the middle by small mirrors. The beams then recombine on the far side of the splitter in an eyepiece, producing a pattern of interference fringes whose transverse displacement depends on the difference in time it takes light to transit the longitudinal and the transverse arms of the apparatus. By analogy with Equation (3.9), this time difference is

$$\Delta t = \frac{l_0 v_e^2}{c^3}, \quad (3.59)$$

where l_0 is the length of the arms, v_e the orbital velocity of the Earth, and c the velocity of light in vacuum. (Henceforth, c refers exclusively to the velocity of light in vacuum.) By turning their apparatus through 90° , the experimentalists expected to reverse the time difference, and, thus, to generate a shift in the interference fringes. The magnitude of this shift is $2c\Delta t/\lambda$ fringes, where $\lambda = 5 \times 10^{-7}$ m is a typical wavelength of light. Thus, given that the lengths of the arms in the experiment were (effectively) 10 m, the expected shift is 0.39 fringes. Such a shift should have been easily measurable. However, no such shift was observed. One possible explanation for this null result is that the Earth completely drags the aether in its immediate vicinity. However, this explanation is in conflict with Bradley and Airy's observations of stellar aberration, as well as Fizeau's experiment.

3.1.6 Lorentz-Fitzgerald Contraction

In 1889, George FitzGerald, followed by Hendrik Lorentz in 1892, suggested that an object moving with speed v with respect to the aether suffers a contraction in length by a factor $\sqrt{1 - v^2/c^2}$ in the direction parallel to the motion, but suffer no contraction in the perpendicular directions. This so-called *length contraction* hypothesis explains the null result in the Michelson-Morley experiment. To be more exact, and referring to Figure 3.2, when the light traverses the leg of the apparatus that is parallel to its velocity with respect to the aether then it takes a time

$$t_1 = \frac{l_0 \sqrt{1 - v_e^2/c^2}}{c - v_e} + \frac{l_0 \sqrt{1 - v_e^2/c^2}}{c + v_e} = \frac{2l_0}{c} \frac{1}{\sqrt{1 - v_e^2/c^2}}, \quad (3.60)$$

where l_0 is the uncontracted length of the leg. On the other hand, when the light traverses the leg of the apparatus that is perpendicular to its velocity with respect to the aether then it takes a time

$$t_2 = \frac{2l_0}{\sqrt{c^2 - v_e^2}} = \frac{2l_0}{c} \frac{1}{\sqrt{1 - v_e^2/c^2}}. \quad (3.61)$$

It can be seen that $t_1 = t_2$, which explains the null result of the Michelson-Morley experiment.

3.1.7 Kennedy-Thorndike Experiment

Suppose that we were to perform a version of the Michelson-Morley experiment in which the two legs of the apparatus are of unequal (uncontracted) lengths l_1 and l_2 . Taking length contraction into account, the time required for light to traverse the first leg of the apparatus is

$$t_1 = \frac{2l_1}{c} \frac{1}{\sqrt{1 - v^2/c^2}}, \quad (3.62)$$

where v is the speed of the laboratory with respect to the aether rest frame. Likewise, the time required for light to traverse the second leg of the apparatus is

$$t_2 = \frac{2l_2}{c} \frac{1}{\sqrt{1 - v^2/c^2}}. \quad (3.63)$$

Hence, the difference between these two times is

$$t_2 - t_1 = \frac{2(l_2 - l_1)}{c} \frac{1}{\sqrt{1 - v^2/c^2}} \approx \frac{2(l_2 - l_1)}{c} + \frac{l_2 - l_1}{c} \frac{v^2}{c^2}. \quad (3.64)$$

Note that the time difference depends on v . Suppose that the laboratory is located on the Earth's equator. In this case, the actual speed of the laboratory with respect to the aether rest frame varies from $v = v_e - v_\Omega$ to $v = v_e + v_\Omega$, throughout the course of a day, where $v_e = 2.977 \times 10^4 \text{ m s}^{-1}$ is the Earth's mean orbital velocity, specified in Equation (3.20), whereas

$$v_\Omega = \Omega R_e = 4.650 \times 10^2 \text{ m s}^{-1} \quad (3.65)$$

is the speed of the Earth's surface due its axial rotation. Here, $\Omega = 7.292 \times 10^{-5} \text{ rad s}^{-1}$ is the Earth's diurnal angular velocity [see Equation (1.351)], and $R_e = 6.378 \times 10^3 \text{ m}$ its equatorial radius. Thus, v varies by about 3% during the course of the day. This variation leads to a variation in the time difference, (3.64), that should be easily measurable. However, when Roy Kennedy and Edward Thorndike performed this experiment in 1932 they observed no variation in the time difference.

3.2 Theoretical Basis of Special Relativity

3.2.1 Postulates of Special Relativity

We have seen that experimental observations of stellar aberration and the speed of light in moving refractive media, together with the results of the Michelson-Morley and the Kennedy-Thorndike experiments, cannot be reconciled within a theoretical framework in which light is assumed to propagate at a fixed speed with respect to an aether.

In 1905, Albert Einstein reconciled all of the aforementioned results within a new theoretical framework known as *special relativity*. The two postulates of this framework are:

1. The laws of physics are invariant (i.e., take equivalent forms) in all inertial frames of reference.
2. The speed of light in vacuum is the same in all inertial frames of reference, irrespective of the motion of the source or the receiver.

Postulate 1 is motivated by the observation that Newton's laws of motion take equivalent forms in all inertial reference frames. (See Section 1.5.4.) One corollary of this observation is that no identical experiment in Newtonian dynamics, performed in various different inertial frames, can provided a way to distinguish one frame from another. Einstein simply generalized this equivalence principle by assuming that it applies to all laws of physics. Thus, according to Einstein, all laws of physics take equivalent forms in all inertial reference frames, implying that no identical experiment, of any kind, performed in various different inertial frames, can provided a way to distinguish one frame from another. Incidentally, the nomenclature 'relativity' derives from the fact that it is impossible to determine that any given inertial frame constitutes an absolute standard of rest; in this respect, all motion is relative.

Postulate 2 follows from Einstein's rejection of the idea of an aether. Einstein assumed that an electromagnetic wave is a self-perpetuating disturbance of electric and magnetic fields that is capable of propagating through a vacuum without the need for a medium. Suppose that we measure the speed of light in vacuum in various different inertial reference frames. If the results of these identical experiments give different speeds then we have found a way of distinguishing the various reference frames from one another. In fact, we could provide a distinct label for each reference frame in terms of its associated speed of light in vacuum. However, this state of affairs is forbidden by Einstein's first postulate. Hence, the speed of light in vacuum must be the same in all of the reference frames.

The speed of light cannot depend on the motion of the source, because, if it did, then we could place a stationary source in each possible inertial reference frame, and then distinguish different frames from one another on the basis of the different speeds of the light emitted by these sources, and measured by a stationary receiver in a particular reference frame. However, this state of affairs is forbidden by Einstein's first postulate. Furthermore, the speed of light cannot depend on the motion of the receiver, because, if it did, then we could place a stationary source in a particular inertial reference frame, and then distinguish different inertial frames from one another on the basis of the different speeds of the light emitted by this source, and measured by stationary receivers in the latter frames.

Incidentally, if light waves propagate through a vacuum with the same speed *in all directions* in one inertial reference frame then they must do so in all inertial reference frames, otherwise we could distinguish the former reference frame from the others, which is contrary to Einstein's first postulate. However, Maxwell's equations predict that light waves propagate with the same speed in all directions in the (presumably inertial) frame of reference in which they are formulated. (See Section 2.4.4.) Hence, a more precise version of Einstein's second postulate is that light waves propagate with an invariant speed in all directions in all inertial frames of reference, irrespective of the motion of the source or the receiver.

Let us consider whether Einstein's first postulate also demands that the speed of sound is the same in all inertial reference frames. Suppose that we measure the speed of sound, in the same gas, in various different inertial reference frames. In general, the measured speed will be different in different reference frames. [See Equation (3.2).] However, this does not violate Einstein's first postulate, because we are not performing the same experiment in the various different reference frames, because each reference frame has a different velocity with respect to the rest frame of the gas. The essential difference between light waves and sound waves is that, because the medium (i.e., aether) with respect to which light waves in a vacuum propagate does not exist, this non-existent medium does not have an identifiable rest frame, whereas the gas through which sound waves propagate always has an identifiable rest frame.

Incidentally, because the Michelson-Morley and Kennedy-Thorndike experiments are, in effect, trying to measure the difference between the speeds of light in vacuum in two different inertial frames, Einstein's second postulate guarantees that they should both give null results.

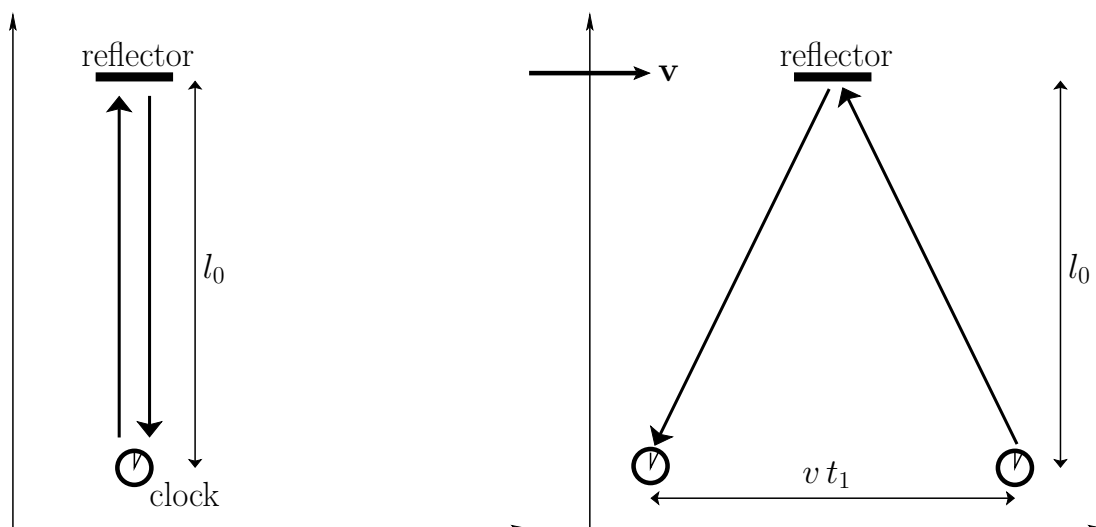


Figure 3.6: Time dilation.

3.2.2 Invariance of Transverse Lengths

Suppose that we recruit a number of observers, and provide each one with an identical meter stick. Next, let us place the different observers on trolleys that move at different velocities with respect to one another. Each time any given pair of observers make a close approach to one another, suppose that they hold up their meter sticks, orientated such that the sticks are parallel to one another, but perpendicular to their relative velocity, and then make a simultaneous measurement of the relative heights of the two top ends, and the two bottom ends, of the sticks. This is equivalent to a comparison of the lengths of the two sticks. If the two lengths are different then we have found a way of distinguishing one inertial frame from another. In fact, we can provide a unique label for each reference frame in terms of the length of its associated observer's meter stick. However, this state of affairs is forbidden by Einstein's first postulate. Hence, all meter sticks must have the same length. In other words, two observers in two inertial reference frames, moving with respect to one another, will always agree on measurements of lengths orientated perpendicular to their relative motion.

3.2.3 Time Dilation

Consider a clock that is synchronized by bouncing a light ray, that propagates through a vacuum, off a reflector that is located a distance l_0 from the clock. Thus, the time taken for the light ray to travel from the clock to the reflector, and back again,

$$t_0 = \frac{2l_0}{c}, \quad (3.66)$$

corresponds to one tick of the clock.

Suppose that we observe the aforementioned clock in a frame of reference that moves with velocity \mathbf{v} with respect to the clock's rest frame, where the direction of \mathbf{v} is perpendicular to the

path of the light ray in the rest frame. See Figure 3.6. Let t_1 be the time required for a light ray to travel from the clock to the reflector, and back again, in the moving frame. In the moving frame, the clock moves a parallel (to $-\mathbf{v}$) distance $v_1 t$ in this time interval. Note that the transverse distance, l_0 , of the reflector from the clock is the same in both reference frames. (See Section 3.2.2.) It is clear, by symmetry, that in traveling from the clock to the reflector, the light ray in the moving frame has moved a transverse distance l_0 and a parallel (to $-\mathbf{v}$) distance $v t_1/2$. Moreover, the ray travels the same transverse and parallel distances in traveling from the reflector back to the clock. Hence, the net path-length of the light ray is

$$L = 2 \left(\frac{v^2 t_1^2}{4} + l_0^2 \right)^{1/2}. \quad (3.67)$$

Now, because the light ray travels at the speed c in the moving frame, according to Einstein's second postulate, we have

$$t_1 = \frac{L}{c}, \quad (3.68)$$

which implies that

$$t_1 = \frac{2}{c} \left(\frac{v^2 t_1^2}{4} + l_0^2 \right)^{1/2}, \quad (3.69)$$

or

$$t_1^2 = \frac{v^2 t_1^2}{c^2} + \frac{4l_0^2}{c^2} = \frac{v^2 t_1^2}{c^2} + t_0^2, \quad (3.70)$$

where use has been made of Equation (3.66). The previous expression can be rearranged to give

$$t_1 = \frac{t_0}{\sqrt{1 - v^2/c^2}}. \quad (3.71)$$

Thus, we conclude that $t_1 > t_0$. Given that a tick of our clock corresponds to the time required for a light ray to travel from the clock to the reflector, and back again, we conclude that the clock ticks more slowly in the moving reference frame than it does in its rest frame. This effect is known as *time dilation*.

We can also conclude that *any* type of clock, not just a light-clock, will tick more slowly in a moving reference frame than in its rest frame, by the same factor as our light clock, otherwise the same experiment (i.e., measuring the time it takes a light ray to travel a distance $2l_0$ in vacuum using the former type of clock) would produced different results in different inertial frames, which is forbidden by Einstein's first postulate.

Let us define the so-called *Lorentz factor*,

$$\gamma = \frac{1}{\sqrt{1 - v^2/c^2}}. \quad (3.72)$$

Note that $\gamma \geq 1$. The time dilation law, (3.71), can be written

$$t_1 = \gamma t_0. \quad (3.73)$$

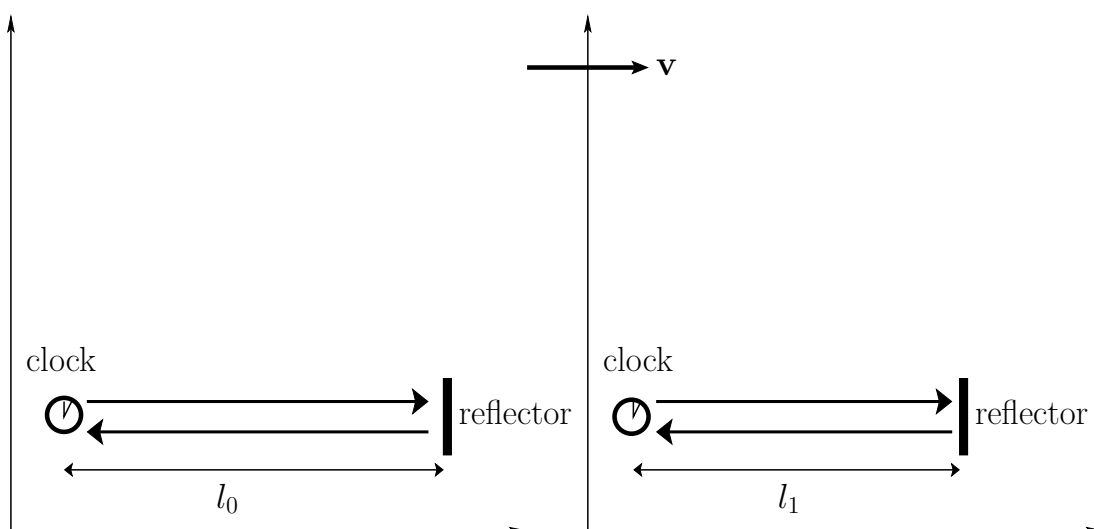


Figure 3.7: Length contraction.

In other words, time is dilated by the Lorentz factor in a moving frame of reference.

Muons are unstable particles that have measured lifetimes of $2.917 \mu\text{s}$ in their rest frame. However, when the Earth's atmosphere is struck by cosmic-ray particles, very energetic muons that move at 98% of the speed of light are produced. The measured lifetimes of these cosmic-ray muons is indeed about five times longer than the rest-frame lifetime of a muon, in accordance with the previous two equations.

3.2.4 Length Contraction

Let us again consider the light-clock introduced in Section 3.2.3. As before, let l_0 be the distance between the clock and the reflector in the clock's rest frame. Thus, each tick of the clock in its rest frame corresponds to the time interval

$$t_0 = \frac{2l_0}{c}. \quad (3.74)$$

[See Equation (3.66).]

Suppose that we observe the aforementioned clock in a frame of reference that moves with velocity \mathbf{v} with respect to the clock's rest frame, where the direction of \mathbf{v} is parallel to the path of the light ray in the rest frame. See Figure 3.7. Let l_1 be the distance between the clock and the reflector in the moving frame. In the moving frame, both the reflector and the clock appear to move at velocity $-\mathbf{v}$ (because they are both at rest in the clock's rest frame, and $-\mathbf{v}$ is the velocity of this frame relative to the moving frame). Suppose that, in the moving frame, the light ray takes a time t_a to travel from the clock to the reflector. The distance traveled by the ray is $l_1 - vt_a$. Hence, because the light ray travels at speed c in the moving frame (irrespective of the motion of the clock or the reflector), according to Einstein's second postulate, we can write

$$t_a = \frac{l_1 - vt_a}{c}, \quad (3.75)$$

which implies that

$$t_a = \frac{l_1}{c + v}. \quad (3.76)$$

Suppose that, in the moving frame, the light ray takes a time t_b to travel from the reflector back to the clock. The distance traveled by the ray is $l + vt_b$. Hence, we can write

$$t_b = \frac{l_1 + vt_b}{c}, \quad (3.77)$$

which implies that

$$t_b = \frac{l_1}{c - v}. \quad (3.78)$$

The net time needed for the light ray to travel from the clock to the reflector, and back again, in the moving frame, is

$$t_1 = t_a + t_b = l_1 \left(\frac{1}{c - v} + \frac{1}{c + v} \right) = \frac{2l_1}{c} \frac{1}{1 - v^2/c^2} = \frac{l_1}{l_0} \frac{t_0}{1 - v^2/c^2}, \quad (3.79)$$

where use has been made of Equation (3.74). This time corresponds to the time interval of the clock's tick in the moving frame. However, we established in Section 3.2.3 that, as a consequence of time dilation,

$$t_1 = \frac{t_0}{\sqrt{1 - v^2/c^2}}. \quad (3.80)$$

Note that if our light clock does not suffer exactly the same time dilation as the clock in Section 3.2.3 then we could distinguish between different inertial frames in terms of the different time dilations suffered by light-clocks in which the light rays traveled parallel and perpendicular to the relative velocities of the frames. However, this state of affairs is prohibited by Einstein's first postulate. The previous two equations yield

$$l_1 = l_0 \sqrt{1 - v^2/c^2} = \frac{l_0}{\gamma}, \quad (3.81)$$

where γ is the Lorentz factor introduced in Equation (3.72). In other words, the distance between the clock and the reflector appears contracted by the Lorentz factor when viewed in the moving frame. Given that we could have placed the clock and the reflector at any two points in space, we conclude that a stationary and a moving observer will not agree on measurements of lengths orientated parallel to their relative motion. In fact, all such lengths will appear contracted by the Lorentz factor to the moving observer. This effect is known as *length contraction*, and is equivalent to the Lorentz-Fitzgerald contraction discussed in Section 3.1.6.

Incidentally, we saw, in Sections 3.1.6 and 3.1.7, that length contraction alone is sufficient to explain the null result of the Michelson-Morley experiment, but not the Kennedy-Thorndike experiment. Hence, we conclude that both length contraction and time dilation are needed to explain the null result of the Kennedy-Thorndike experiment. Another way of saying this is that the null result of the Michelson-Morley experiment can be regarded as experimental verification of length contraction, whereas the null result of the Kennedy-Thorndike experiment can be regarded as experimental validation of time dilation. It should be noted that both of these experiments have been repeated many times, over the years, and that the null results of the experiments are now established to very great accuracy.

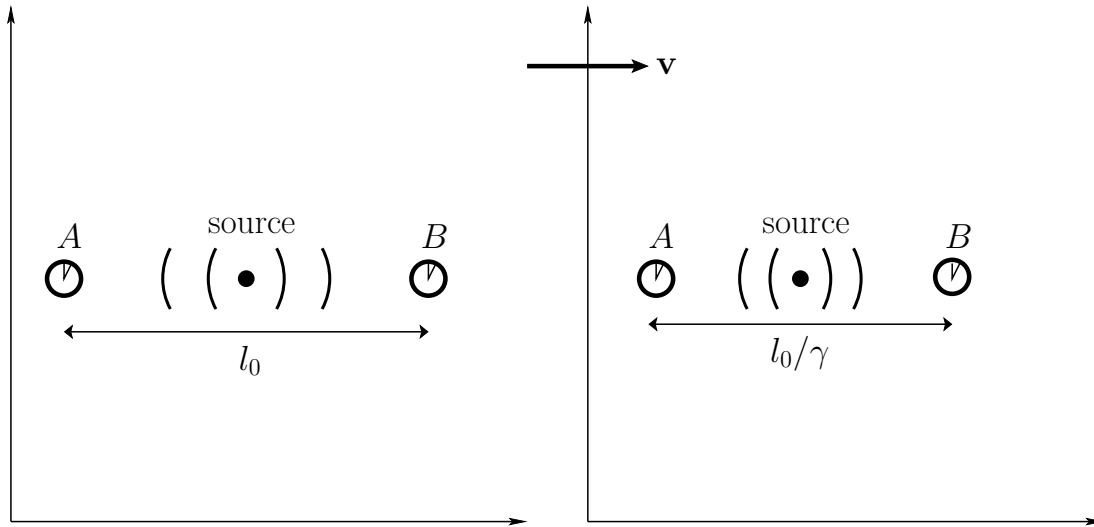


Figure 3.8: Clock error.

3.2.5 Clock Error

Consider two clocks, A and B , that are located a distance l_0 apart in their mutual rest frame. Suppose that the clocks are synchronized using light pulses emitted from a source that lies half-way between them. Let us observe the clocks in a reference frame that moves with velocity \mathbf{v} with respect to the clocks' rest frame in a direction that is parallel to their mutual displacement. See Figure 3.8. In the moving frame, the contracted distance between the two clocks is l_0/γ , but the source is still located half-way between the clocks. Moreover, the two clocks appear to move with the same velocity, $-\mathbf{v}$. Consider a light pulse that is emitted by the source and travels to the two clocks. Suppose that, in the moving frame, it takes a time t_a for the pulse to travel from the source to clock A . The pulse travels a distance $l_0/(2\gamma) + vt_a$. Thus, given that the pulse travels at the speed c , according to Einstein's second postulate, we have

$$t_a = \frac{l_0/(2\gamma) + vt_a}{c}, \quad (3.82)$$

or

$$t_a = \frac{l_0}{2\gamma(c - v)}. \quad (3.83)$$

Suppose that, in the moving frame, it takes a time t_b for the pulse to travel from the source to clock B . The pulse travels a distance $l_0/(2\gamma) - vt_b$. Thus, given that the pulse travels at the speed c , we have

$$t_b = \frac{l_0/(2\gamma) - vt_b}{c}, \quad (3.84)$$

or

$$t_b = \frac{l_0}{2\gamma(c + v)}. \quad (3.85)$$

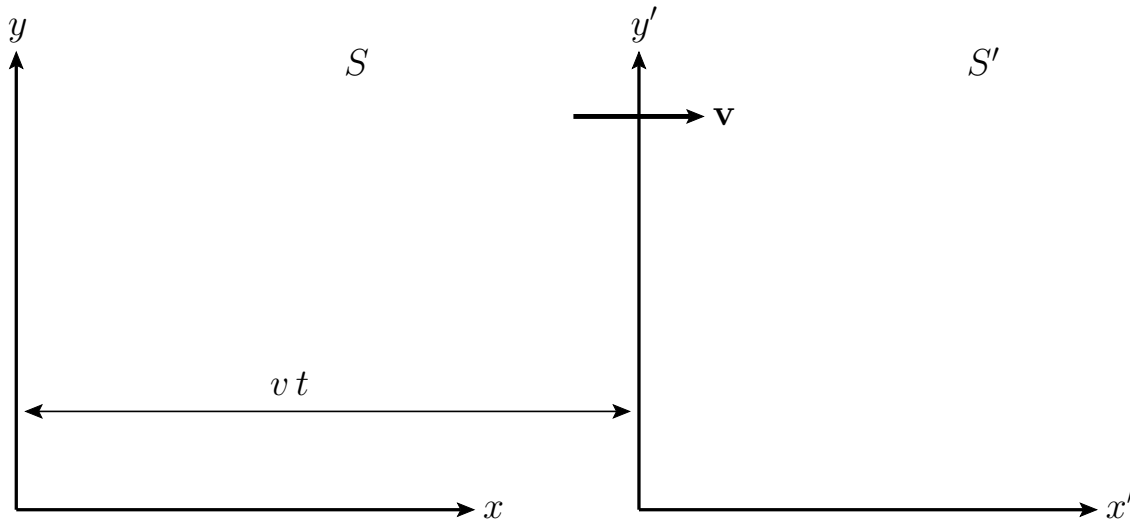


Figure 3.9: Inertial reference frames.

Now, in the clocks' rest frame, the pulse arrives at clocks A and B simultaneously. However, in the moving frame, the pulse arrives at clock B prior to its arrival at clock A (because $t_b < t_a$). In other words, two events, a spatial distance l_0 apart, that take place simultaneously in a particular reference frame, do not appear to take place simultaneously in a reference frame that moves with velocity \mathbf{v} in the direction of the mutual displacement of the two events. This phenomenon is known as *clock error*. The time difference between the two events in the moving frame is

$$\Delta t = t_a - t_b = \frac{l_0}{2\gamma} \left(\frac{1}{c-v} - \frac{1}{c+v} \right) = \frac{l_0}{2\gamma} \frac{2v}{c^2 - v^2}. \quad (3.86)$$

which reduces to

$$\Delta t = \frac{\gamma v l_0}{c^2}. \quad (3.87)$$

3.2.6 Galilean Transformation

Consider two inertial frames of reference, S and S' . Let frame S' move at velocity \mathbf{v} with respect to frame S . Let us set up right-handed Cartesian coordinate systems in both frames. Suppose that the coordinate systems are in the so-called *standard configuration* in which the corresponding coordinate axes are parallel, the x -axis in each system is parallel to \mathbf{v} , and the origins of the systems coincide at time $t = 0$. See Figure 3.9. Consider an instantaneous 'event' with a definite spatial location, such as the flashing of a light-bulb. Suppose that the event occurs at time t and has displacement (x, y, z) in frame S . Suppose that the event occurs at time t' and has displacement (x', y', z') in frame S' . What is the relationship between (x, y, z, t) and (x', y', z', t') . Well, according to standard Newtonian physics, the "common sense" relationship between the two sets of coordinates is

$$x' = x - vt, \quad (3.88)$$

$$y' = y, \quad (3.89)$$

$$z' = z, \quad (3.90)$$

$$t' = t. \quad (3.91)$$

(See Section 1.5.4.) As we have already mentioned, this transformation of coordinates is known as the *Galilean transformation*. Consider, now, a moving event whose coordinates in S are $x = x(t)$, $y = y(t)$, and $z = z(t)$. The Cartesian components of the instantaneous velocity of our event in S are $u_x = dx/dt$, $u_y = dy/dt$, $u_z = dz/dt$, whereas the corresponding components in S' are $u'_x = dx'/dt'$, $u'_y = dy'/dt'$, $u'_z = dz'/dt'$. Hence, we can derive the following Galilean velocity transformation from Equations (3.88)–(3.91):

$$u'_x = u_x - v, \quad (3.92)$$

$$u'_y = u_y, \quad (3.93)$$

$$u'_z = u_z. \quad (3.94)$$

However, if the event in question is the path of a light ray that moves with velocity $c \mathbf{e}_x$ in S then, according to the Galilean velocity transform, the light ray moves with velocity $(c - v) \mathbf{e}_x$ in S' . In other words, the light ray travels at different speeds in the two frames of reference. However, this state of affairs is forbidden by Einstein's first postulate. Hence, we deduce that the Galilean transformation, (3.88)–(3.91), is actually inconsistent with the theory of relativity.

3.2.7 Lorentz Transformation

Let us see if we can derive a transformation of coordinates that is consistent with relativity. In frame S , at time t , the origin of the x' -axis is located a perpendicular distance vt from the origin of the x -axis. Moreover, in S , our event is located a perpendicular distance x'/γ from the origin of the x' -axis, because of length contraction. [See Equation (3.81).] Thus,

$$x = vt + \frac{x'}{\gamma}, \quad (3.95)$$

which yields

$$x' = \gamma(x - vt). \quad (3.96)$$

Because there is no motion between the two frames in the y -direction or the z -direction, and because there is no contraction of lengths perpendicular to the x -axis, we deduce that

$$y' = y, \quad (3.97)$$

$$z' = z. \quad (3.98)$$

Finally, in S' , a clock located at the origin reads γt when a clock at the origin of S reads t , as a consequence of time dilation (note that, in S' , a clock in S appears to run slowly by a factor γ). [See Equation (3.73).] Furthermore, a second clock in S' , displaced from the origin a distance x

(measured in S) in the x -direction, appears to read $\gamma t - \gamma v x/c^2$ as a consequence of clock error. [See Equation (3.87).] Hence,

$$t' = \gamma t - \frac{\gamma v x}{c^2} = \gamma \left(t - \frac{v x}{c^2} \right). \quad (3.99)$$

We deduce that the transformation of coordinates that is consistent with the theory of relativity is

$$x' = \gamma(x - vt), \quad (3.100)$$

$$y' = y, \quad (3.101)$$

$$z' = z, \quad (3.102)$$

$$t' = \gamma \left(t - \frac{v x}{c^2} \right). \quad (3.103)$$

This transformation is known as the *Lorentz transformation*. Note that, in the limit in which the relative velocity of frames S and S' is *non-relativistic* (i.e., much smaller in magnitude than the speed of light in vacuum), so that $v/c \ll 1$, and $\gamma \rightarrow 1$, the Lorentz transformation morphs into the Galilean transformation, (3.88)–(3.91). Thus, the common sense transformation, (3.88)–(3.91), holds as long as the relative velocity between the two frames of reference is much smaller than the velocity of light in vacuum. Incidentally, it is easily shown from Equations (3.100)–(3.103) that

$$x = \gamma(x' + vt'), \quad (3.104)$$

$$y = y', \quad (3.105)$$

$$z = z', \quad (3.106)$$

$$t = \gamma \left(t' + \frac{v x'}{c^2} \right). \quad (3.107)$$

In the Galilean transformation, (3.88)–(3.91), the transformation of time is completely independent from that of space. This is no longer the case in the Lorentz transformation, (3.100)–(3.103). In fact, the transformations of space and time are mixed together in special relativity in such a manner that, rather than thinking of space and time as separate concepts, it makes more sense to talk about a generalized concept that Einstein called *spacetime*.

3.2.8 Spacetime Interval

Consider two events, 1 and 2, whose spacetime coordinates in some inertial frame S are (x_1, y_1, z_1, t_1) and (x_2, y_2, z_2, t_2) , respectively. Let us form the differences between these coordinates, $\Delta x = x_2 - x_1$, $\Delta y = y_2 - y_1$, $\Delta z = z_2 - z_1$, and $\Delta t = t_2 - t_1$. The spatial distance, Δd , between the two events is written

$$(\Delta d)^2 = (\Delta x)^2 + (\Delta y)^2 + (\Delta z)^2. \quad (3.108)$$

Suppose that we shift the origin of our coordinate system in S . It is obvious that this process does not change the distance between our two events. In other words, Δd is invariant under a shift

of the origin of the coordinate system, as is easily verified. Suppose that we rotate our coordinate axes in S . (See Section A.5.) Such a process is length-preserving. In other words, Δd is invariant under a rotation of the coordinate axes, as is easily verified. However, it is evident, by inspection, that Δd is not invariant under a Lorentz transformation. Let us try to find a quantity that is invariant.

Consider a second inertial frame, S' , that moves with velocity $\mathbf{v} = v \mathbf{e}_x$ with respect to S , and is also in a standard configuration with respect to S . Let events 1 and 2 have spacetime coordinates (x'_1, y'_1, z'_1, t'_1) and (x'_2, y'_2, z'_2, t'_2) , respectively, in S' . Let us again form the differences between these coordinates, $\Delta x' = x'_2 - x'_1$, $\Delta y' = y'_2 - y'_1$, $\Delta z' = z'_2 - z'_1$, and $\Delta t' = t'_2 - t'_1$. The *spacetime interval* between events 1 and 2 in S is defined

$$(\Delta s)^2 = c^2 (\Delta t)^2 - (\Delta x)^2 - (\Delta y)^2 - (\Delta z)^2. \quad (3.109)$$

The Lorentz transformation, (3.100)–(3.103), yields

$$\Delta x' = \gamma (\Delta x - v \Delta t), \quad (3.110)$$

$$\Delta y' = \Delta y, \quad (3.111)$$

$$\Delta z' = \Delta z, \quad (3.112)$$

$$\Delta t' = \gamma \left(\Delta t - \frac{v \Delta x}{c^2} \right). \quad (3.113)$$

Hence, the spacetime interval between the two events in S' is

$$\begin{aligned} (\Delta s')^2 &= c^2 (\Delta t')^2 - (\Delta x')^2 - (\Delta y')^2 - (\Delta z')^2 \\ &= \gamma^2 \left[c^2 (\Delta t)^2 - 2v \Delta x \Delta t + \frac{v^2}{c^2} (\Delta x)^2 \right] - \gamma^2 [(\Delta x)^2 - 2v \Delta x \Delta t + v^2 (\Delta t)^2] \\ &\quad - (\Delta y)^2 - (\Delta z)^2 \\ &= c^2 (\Delta t)^2 \gamma^2 \left(1 - \frac{v^2}{c^2} \right) - (\Delta x)^2 \gamma^2 \left(1 - \frac{v^2}{c^2} \right) - (\Delta y)^2 - (\Delta z)^2 \\ &= c^2 (\Delta t)^2 - (\Delta x)^2 - (\Delta y)^2 - (\Delta z)^2 \\ &= (\Delta s)^2, \end{aligned} \quad (3.114)$$

where use has been made of Equation (3.72). Thus, we conclude that the spacetime interval between two events is invariant under a Lorentz transformation. In other words, the interval is the same in all inertial reference frames. What is the significance of this result? Suppose that a light ray travels in a straight-line from event 1 to event 2. The speed of the light ray in S is

$$\frac{\Delta d}{\Delta t} = c \left[1 - \frac{(\Delta s)^2}{(c \Delta t)^2} \right]^{1/2}, \quad (3.115)$$

where use has been made of Equations (3.108) and (3.109). However, we know, from Einstein's second postulate, that this speed is equal to c . Hence, we deduce that

$$\Delta s = 0. \quad (3.116)$$

In other words, if a light ray travels between two events in spacetime then the interval between these two events is zero. This implies that the interval between the two events is zero in all inertial frames of reference. In particular, the interval in S' is $\Delta s' = 0$. Now, the speed of the light ray in S' is

$$\frac{\Delta d'}{\Delta t'} = c \left[1 - \frac{(\Delta s')^2}{(c \Delta t')^2} \right]^{1/2} = c. \quad (3.117)$$

Hence, we deduce that the invariance of the spacetime interval under Lorentz transformation guarantees that light travels through a vacuum at the speed c in all inertial frames of reference.

3.2.9 Transformation of Velocity

Consider a particle in some inertial reference frame, S , that moves at the fixed (subluminal) velocity $\mathbf{u} = (u_x, u_y, u_z)$. Suppose that the particle is located at the origin at time $t = 0$. It follows that, at time t , the particle is located at point $\mathbf{r} = (u_x t, u_y t, u_z t)$. Let us observe the motion of the particle in a second inertial frame, S' , that moves with (subluminal) velocity $\mathbf{v} = v \mathbf{e}_x$ with respect to S , and is in a standard configuration with respect to S . It follows from Equations (3.100)–(3.103) that the particle is located at the origin of S' at $t' = 0$. The location of the particle in S' , at time t' , can be written $\mathbf{r}' = (u'_x t', u'_y t', u'_z t')$, where $\mathbf{u}' = (u'_x, u'_y, u'_z)$ is the particle's velocity in S' . It follows from Equations (3.100)–(3.103) that

$$u'_x t' = \gamma(u_x t - v t), \quad (3.118)$$

$$u'_y t' = u_y t, \quad (3.119)$$

$$u'_z t' = u_z t, \quad (3.120)$$

$$t' = \gamma \left(t - \frac{v u_x t}{c^2} \right), \quad (3.121)$$

which yields

$$u'_x = \frac{u_x - v}{1 - u_x v/c^2}, \quad (3.122)$$

$$u'_y = \frac{u_y}{\gamma(1 - u_x v/c^2)}, \quad (3.123)$$

$$u'_z = \frac{u_z}{\gamma(1 - u_x v/c^2)}. \quad (3.124)$$

This result is known as the *transformation of velocity*.

Let $u = (u_x^2 + u_y^2 + u_z^2)^{1/2}$ and $u' = (u'^2_x + u'^2_y + u'^2_z)^{1/2}$ be the speeds of the particle in frames S and S' , respectively. It is easily demonstrated, from the transformation of velocity, that

$$c^2 - u'^2 = \frac{c^2(c^2 - u^2)(c^2 - v^2)}{(c^2 + u_x v)^2}. \quad (3.125)$$

If $|u| < c$ and $|v| < c$ then the right-hand side is positive, implying that $|u'| < c$. In other words, the resultant of two subluminal velocities is another subluminal velocity. It is evident that a particle can

never attain the velocity of light relative to a given inertial frame, no matter how many subluminal velocity increments it is given. It follows that no inertial frame can ever appear to propagate with a superluminal velocity with respect to any other inertial frame (because we can track a given inertial frame in terms of a particle that remains at rest at the origin of that frame). Of course, if $|u| = c$ then $|u'| = c$. In other words, a particle traveling at the speed of light in one inertial frame does so in all inertial frames.

It is evident from Equation (3.125) that there is only a *single* speed—namely, $u = c$ —that is the same in all inertial frames of reference. Now, according to Einstein's first postulate, *any* wave that propagates in the absence of a physical medium must propagate at the same speed in all inertial frames of reference, otherwise the different wave speeds in different reference frames could be used to distinguish between the frames. Hence, we deduce that all waves that propagate in the absence of a physical medium (e.g., a gas, liquid, or solid) must propagate at the common speed c in all inertial reference frames. Thus, gravitational waves, which are ripples in the fabric of spacetime, must travel at the same speed, c , as electromagnetic waves, because both waves propagate in the absence of media. Thus, we could just as well designate c as the speed of gravitational waves.

Note, finally, that the Lorentz transformation is the only (linear) transformation of coordinates that preserves the speed c , and morphs into the tried and tested Galilean transformation in the limit that $v/c \ll 1$. In fact, it is possible to guess the form of the Lorentz transformation by searching for a (linear) coordinate transformation that has these two properties.

3.2.10 Causality

Let events 1 and 2 have spacetime coordinates $(x_1, 0, 0, t_1)$ and $(x_2, 0, 0, t_2)$ in some inertial reference frame, S . Suppose that event 1 causes event 2. It follows that $t_1 < t_2$. In other words, event 1 necessarily precedes event 2 in time. Let

$$u = \frac{x_2 - x_1}{t_2 - t_1} \quad (3.126)$$

be the velocity with which information flows from event 1 to event 2 in order to allow the former event to cause the latter. Let us observe the two events in a second inertial frame, S' , that moves at velocity $\mathbf{v} = v \mathbf{e}_x$ with respect to S , and is in a standard configuration with respect to S . According to Equation (3.103),

$$t'_2 - t'_1 = \gamma \left(t_2 - \frac{v x_2}{c^2} \right) - \gamma \left(t_1 - \frac{v x_1}{c^2} \right), \quad (3.127)$$

or

$$t'_2 - t'_1 = \gamma (t_2 - t_1) \left(1 - \frac{u v}{c^2} \right). \quad (3.128)$$

Now, irrespective of the value of v , whose magnitude can never exceed c , event 2 can never occur prior to event 1 in S' , otherwise we could classify inertial frames into two groups; those in which event 1 appears to cause event 2, and those in which event 2 appears to cause event 1. However, this state of affairs is forbidden by Einstein's first postulate. Thus, we require $t'_2 - t'_1 > 0$ for all $|v| < c$. It is clear from Equation (3.128) that this is only possible if $|u| < c$. Hence, we deduce that information can never propagate faster than the speed of light in vacuum, in any inertial reference

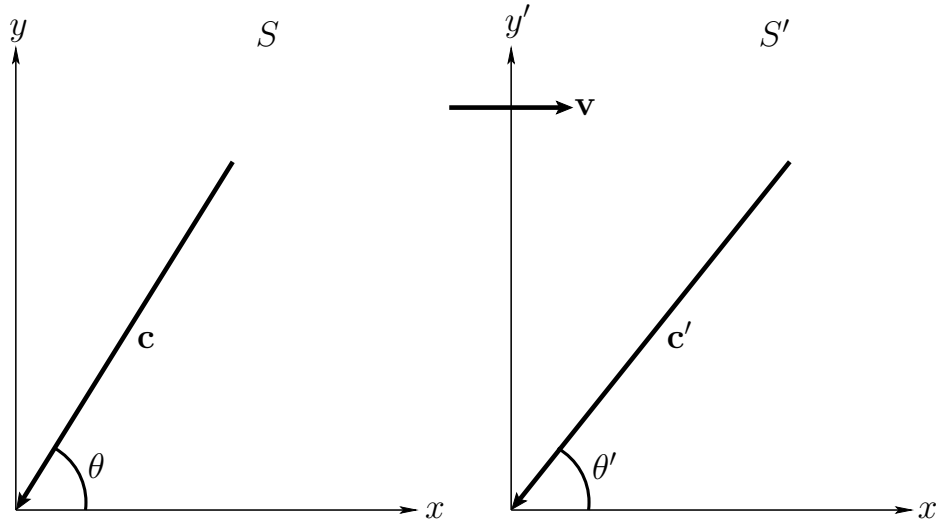


Figure 3.10: Relativistic aberration of light.

frame, otherwise it is possible to find other inertial reference frames in which causality appears to be violated.

3.2.11 Relativistic Aberration of Light

Consider a light ray that travels from a distant source to an observer located at the origin of some inertial frame, S . Let \mathbf{c} be the phase velocity of the light ray. Of course, $|\mathbf{c}| = c$, where c is the speed of light in vacuum. Suppose that \mathbf{c} lies in the x - y plane, such that its direction subtends an angle θ with the $-x$ -direction, as shown in Figure 3.10. It is clear from the figure that $c_x = -c \cos \theta$ and $c_y = -c \sin \theta$. Suppose that a second observer, moving with velocity $\mathbf{v} = v \mathbf{e}_x$ with respect to the first, observes the light ray. Let \mathbf{c}' be the phase velocity of the light ray in the second observer's frame, S' , which is in a standard configuration with respect to frame S . Of course, $|\mathbf{c}'| = c$. Suppose that \mathbf{c}' lies in the x' - y' plane, such that its direction subtends an angle θ' with the $-x'$ -direction, as shown in Figure 3.10. It is clear from the figure that $c'_x = -c \cos \theta'$ and $c'_y = -c \sin \theta'$. The transformation of velocity, (3.122)–(3.124), yields

$$\tan \theta' = \frac{-u'_y}{-u'_x} = \frac{-u_y}{-\gamma(u_x - v)} = \frac{c \sin \theta}{-\gamma(-c \cos \theta - v)}, \quad (3.129)$$

or

$$\tan \theta' = \frac{\sin \theta}{\gamma(\cos \theta + v/c)}. \quad (3.130)$$

Thus, the direction of the light ray, and, hence the angular position of the source, appears different to the two observers.

In particular, suppose that the first observer is located in the rest frame of the Sun, and the second is located on the Earth, whose instantaneous orbital velocity about the Sun is $\mathbf{v} = v_e \mathbf{e}_x$,

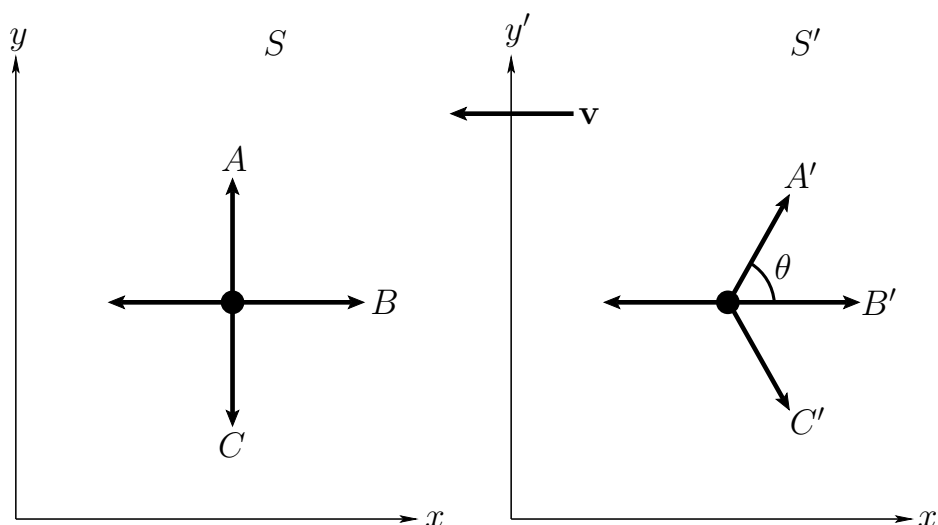


Figure 3.11: Relativistic beaming of light.

where $v_e = 2.977 \times 10^4 \text{ m s}^{-1}$. In this case, the previous equation yields

$$\tan \theta' = \frac{\sin \theta}{\gamma (\cos \theta + \kappa)}, \quad (3.131)$$

where

$$\kappa = \frac{v_e}{c} = 9.930 \times 10^{-5}, \quad (3.132)$$

$$\gamma = \frac{1}{\sqrt{1 - v_e^2/c^2}} = 1.000000005. \quad (3.133)$$

It can be seen that formula (3.131) is almost indistinguishable from the classical aberration formula, (3.23). Thus, it is clear that special relativity is capable of accounting for Bradley's observation of the aberration of starlight. (See Section 3.1.3.) Furthermore, it is obvious that the relativistic aberration of light is associated with the properties of the Lorentz transformation between two frames of reference, moving with respect to one another, rather than the velocity of light with respect to any particular medium. It follows that special relativity is also capable of accounting for Airy's observation of the aberration of starlight. (See Section 3.1.4.)

3.2.12 Relativistic Beaming of Light

Consider a point source that emits light isotropically in all directions in its rest frame, S . Let us observe this source in a frame of reference, S' , that moves with velocity $\mathbf{v} = -v \mathbf{e}_x$, and is in a standard configuration, with respect to frame S . Thus, the source appears to move with velocity $v \mathbf{e}_x$ in frame S' . See Figure 3.11. Now, half of the emitted radiation in S emerges in the region ABC , bounded by the rays A and C shown in the figure. Likewise, half the emitted radiation in

S' emerges in the region $A'B'C'$, bounded by the rays A' and C' shown in the figure. Ray A has the phase velocity $(0, c, 0)$. Likewise, ray A' has the phase velocity $(c \cos \theta, c \sin \theta, 0)$, where the angle θ is shown in the figure. By symmetry, the angle subtended between A' and B' is the same as that subtended between C' and B' . The transformation of velocity, (3.122)–(3.124), yields

$$c \cos \theta = v, \quad (3.134)$$

$$c \sin \theta = \frac{c}{\gamma}, \quad (3.135)$$

or

$$\sin \theta = \frac{1}{\gamma}. \quad (3.136)$$

It follows, that in a frame of reference in which the source moves with velocity \mathbf{v} , half of the emitted radiation is beamed into a cone whose axis is \mathbf{v} , and whose half-angle is $\sin^{-1}(1/\gamma)$. If the source is moving very close to the velocity of light then $\gamma \gg 1$, and $\theta \simeq 1/\gamma \ll 1$. In other words, the emitted radiation is beamed very strongly in the direction of motion of the source.

3.2.13 Light Propagation through Dielectric Media

Consider a transparent dielectric medium, such as air, water, or glass. Let n be the refractive index of the medium. Thus, in the rest frame, S , of the medium, light propagates at the phase velocity c/n . Let us transform to a frame of reference, S' , that moves at velocity $\mp v \mathbf{e}_x$, and is in a standard configuration, with respect to frame S . Thus, the medium appears to flow at the velocity $\pm v \mathbf{e}_x$ in frame S' . According to Equation (3.122), the phase velocity of light propagating in the $+x$ -direction in frame S' , in which the medium flows at velocity $\pm v \mathbf{e}_x$, is

$$u_{\pm} = \frac{c/n \pm v}{1 \pm v/(cn)} = \frac{c}{n} \pm \frac{v(1 - 1/n^2)}{1 \pm v/(cn)}. \quad (3.137)$$

In the limit in which the flow velocity of the medium is much smaller than the velocity of light in vacuum, $v/c \ll 1$, the previous equation reduces to

$$u_{\pm} \simeq \frac{c}{n} \pm v \left(1 - \frac{1}{n^2}\right). \quad (3.138)$$

However, this is identical to the phase velocity in dielectric media measured by Fizeau. (See Section 3.1.4.) Thus, we can now appreciate that special relativity is capable of accounting for all of the experimental observations discussed in Section 3.1.

3.3 Relativistic Dynamics

3.3.1 Transformation of Acceleration

Consider a particle that moves with constant acceleration $\mathbf{a}' = a_0 \mathbf{e}_x$ in its instantaneous rest frame, S' . Let the particle be located at the origin of S' at time $t' = 0$. It follows that a very short time, δt ,

later, the particle's spacetime coordinates in S' are

$$x' = \frac{1}{2} a_0 \delta t^2, \quad (3.139)$$

$$y' = 0, \quad (3.140)$$

$$z' = 0, \quad (3.141)$$

$$t' = \delta t. \quad (3.142)$$

Consider a second frame of reference, S , that moves with velocity $-v \mathbf{e}_x$ with respect to S' , and is also in a standard configuration with respect to S' . Thus, the particle's instantaneous velocity in S is $\mathbf{v} = v \mathbf{e}_x$. In S , the particle moves from the origin at time $t = 0$, to a point whose spacetime coordinates are

$$x = \gamma(x' + vt') = \gamma \left(\frac{1}{2} a_0 \delta t^2 + v \delta t \right), \quad (3.143)$$

$$y = 0, \quad (3.144)$$

$$z = 0, \quad (3.145)$$

$$t = \gamma \left(t' + \frac{v x'}{c^2} \right) = \gamma \left(\delta t + \frac{v a_0 \delta t^2}{2 c^2} \right), \quad (3.146)$$

where $\gamma = (1 - v^2/c^2)^{-1/2}$, a very short time later. Here, use has been made of Equations (3.104)–(3.107), as well as Equations (3.139)–(3.142). If $\mathbf{a} = a \mathbf{e}_x$ is the particle's instantaneous acceleration in S' then we expect the relation

$$x = vt + \frac{1}{2} a t^2 \quad (3.147)$$

to hold for a short time interval (i.e., in the limit $\delta t \rightarrow 0$). It follows from Equations (3.143) and (3.146) that

$$\gamma \left(\frac{1}{2} a_0 \delta t^2 + v \delta t \right) = v \gamma \left(\delta t + \frac{v a_0 \delta t^2}{2 c^2} \right) + \frac{1}{2} a \gamma^2 \left(\delta t + \frac{v a_0 \delta t^2}{2 c^2} \right)^2. \quad (3.148)$$

Note that the terms in the previous equation that are first order in δt cancel one another. Equating the terms that are second order in δt , we obtain

$$\gamma a_0 = \gamma \frac{v^2}{c^2} a_0 + \gamma^2 a, \quad (3.149)$$

or

$$a = \frac{a_0 (1 - v^2/c^2)}{\gamma}, \quad (3.150)$$

which reduces to

$$a = \frac{a_0}{\gamma^3}. \quad (3.151)$$

Thus, we conclude that the particle's instantaneous acceleration in a frame of reference in which it has a finite speed is always less than that in its instantaneous rest frame.

Given that, by definition,

$$\frac{dv}{dt} = a, \quad (3.152)$$

the particle's equation of motion in S , which is assumed to be an inertial frame, is

$$\frac{dv}{dt} = \frac{a_0}{\gamma^3} = a_0 \left(1 - \frac{v^2}{c^2}\right)^{3/2}. \quad (3.153)$$

If $v = 0$ at $t = 0$ then the previous equation can be integrated to give

$$\frac{v}{(1 - v^2/c^2)^{1/2}} = a_0 t, \quad (3.154)$$

or

$$v = \frac{a_0 t}{(1 + a_0^2 t^2/c^2)^{1/2}}. \quad (3.155)$$

Thus, as seen by an observer at rest in frame S , our particle initially (i.e., for $t \ll c/a_0$) accelerates such that

$$v \simeq a_0 t, \quad (3.156)$$

in accordance with Newtonian dynamics. However, when the speed of the particle becomes comparable with the speed of light in vacuum, the linear increase in speed with time specified in the previous equation breaks down, and, instead,

$$v \rightarrow c \text{ as } t \rightarrow \infty. \quad (3.157)$$

Thus, despite the particle's constant acceleration, a_0 , in its instantaneous rest frame, the particle never appears to move faster than the speed of light to a stationary observer.

3.3.2 Relativistic Equation of Motion

Suppose that the particle discussed in the previous section has a mass m_0 in its instantaneous rest frame. Given that the particle's acceleration in its instantaneous rest frame is $a_0 \mathbf{e}_x$, the particle is clearly subject to a force $\mathbf{f} = f \mathbf{e}_x$, where

$$f = m_0 a_0. \quad (3.158)$$

Thus, according to Equation (3.153), the particle's equation of motion in an inertial reference frame in which its instantaneous velocity is $\mathbf{v} = v \mathbf{e}_x$ is

$$\frac{dv}{dt} = \frac{f}{m_0} \left(1 - \frac{v^2}{c^2}\right)^{3/2} \quad (3.159)$$

However, the previous equation can be rearranged to give

$$f = \frac{d}{dt} \left[\frac{m_0 v}{(1 - v^2/c^2)^{1/2}} \right] = \frac{d(\gamma m_0 v)}{dt}. \quad (3.160)$$

Let us define the *relativistic mass* of the particle as

$$m = \gamma m_0, \quad (3.161)$$

and its *relativistic momentum* as

$$\mathbf{p} = m \mathbf{v}. \quad (3.162)$$

Thus, Equation (3.160) implies that the relativistic equation of motion of the particle is

$$\mathbf{f} = \frac{d\mathbf{p}}{dt}, \quad (3.163)$$

which is analogous in form to Newton's second law of motion, (1.17). Thus, we conclude that the reason that a particle of *rest mass* (i.e., mass in its instantaneous rest frame) m_0 , subject to a constant force \mathbf{f} , never achieves a speed greater than the speed of light is that the particle's relativistic mass, γm_0 , increases as it moves faster, and tends to infinity as its speed approaches the speed of light.

3.3.3 Work and Energy

Suppose that the force $\mathbf{f} = f \mathbf{e}_x$, that acts on the particle discussed in the previous section, causes the particle to displace a distance $d\mathbf{r} = dx \mathbf{e}_x$. The net work done on the particle is clearly

$$dW = \mathbf{f} \cdot d\mathbf{r} = f dx = \frac{d(mv)}{dt} dx = v d(mv), \quad (3.164)$$

because, by definition, $v = dx/dt$. (See Section 1.3.2.) Here, use has been made of Equations (3.160) and (3.161). However,

$$m = \gamma m_0 = \left(1 - \frac{v^2}{c^2} \right)^{-1/2} m_0, \quad (3.165)$$

so

$$v = c \left(1 - \frac{m_0^2}{m^2} \right)^{1/2}. \quad (3.166)$$

The previous equation can be combined with Equation (3.164) to give

$$\begin{aligned} dW &= c^2 \left(1 - \frac{m_0^2}{m^2} \right)^{1/2} d \left[m \left(1 - \frac{m_0^2}{m^2} \right)^{1/2} \right] \\ &= c^2 \left(1 - \frac{m_0^2}{m^2} \right)^{1/2} d \left(\sqrt{m^2 - m_0^2} \right) \end{aligned}$$

$$\begin{aligned}
&= c^2 \left(1 - \frac{m_0^2}{m^2}\right)^{1/2} \frac{m \, dm}{\sqrt{m^2 - m_0^2}} \\
&= c^2 \, dm.
\end{aligned} \tag{3.167}$$

Suppose that the particle is initially at rest, so that its initial relativistic mass is m_0 . Let the force perform net work W on the particle, in the process causing its relativistic mass to increase to m . It is clear from the previous equation that

$$W = (m - m_0) c^2. \tag{3.168}$$

However, we know that the net work that a force does on a particle causes the particle's kinetic energy, K , to increase by a corresponding amount. (See Section 1.3.2.) Thus, given that the particle's initial kinetic energy is zero, we deduce that its kinetic energy is

$$K = (m - m_0) c^2 \tag{3.169}$$

when its relativistic mass is m .

Equation (3.169) can be combined with Equation (3.165) to give

$$K = m_0 c^2 \left[\left(1 - \frac{v^2}{c^2}\right)^{-1/2} - 1 \right]. \tag{3.170}$$

In the limit that the particle is moving at a non-relativistic speed, such that $v/c \ll 1$, the previous equation reduces to

$$K \simeq m_0 c^2 \left[\left(1 + \frac{1}{2} \frac{v^2}{c^2} + \dots\right) - 1 \right], \tag{3.171}$$

or

$$K = \frac{1}{2} m_0 v^2. \tag{3.172}$$

This is consistent with the Newtonian definition of kinetic energy, as long as we identify the rest mass of the particle with its mass in Newtonian dynamics. (See Section 1.3.2.)

3.3.4 Relativistic Energy

Equation (3.169) can be written

$$m c^2 = K + m_0 c^2. \tag{3.173}$$

Let us define the *relativistic energy*, E , of our particle as

$$E = m c^2. \tag{3.174}$$

The previous two equations suggest that the particle possesses two types of energy. First, the particle possesses kinetic energy, K , by virtue of its motion. Second, the particle possesses *rest mass energy*,

$$E_0 = \frac{1}{2} m_0 c^2, \tag{3.175}$$

by virtue of its rest mass. The conjecture that mass is a form of energy was first made by Einstein in 1905. Incidentally, Equation (3.174) implies that conservation of energy in relativistic dynamics is equivalent to conservation of relativistic mass.

3.3.5 Relativistic Energy-Momentum Relation

According to Equations (3.162), (3.165), and (3.174), a particle of rest mass m_0 , moving at velocity \mathbf{v} , has a relativistic momentum

$$\mathbf{p} = \frac{m_0 \mathbf{v}}{\sqrt{1 - v^2/c^2}}, \quad (3.176)$$

and a relativistic energy

$$E = \frac{m_0 c^2}{\sqrt{1 - v^2/c^2}}. \quad (3.177)$$

Thus,

$$\frac{E^2}{c^2} - |\mathbf{p}|^2 = \frac{m_0^2 c^2}{1 - v^2/c^2} - \frac{m_0^2 c^2 (v^2/c^2)}{1 - v^2/c^2} = m_0^2 c^2, \quad (3.178)$$

which leads to the *relativistic energy-momentum relation*,

$$\frac{E^2}{c^2} - |\mathbf{p}|^2 = m_0^2 c^2. \quad (3.179)$$

Now, given that the rest mass is independent of the particle's motion (i.e., it is the same in all inertial frames of reference), we deduce that $E^2/c^2 - |\mathbf{p}|^2$ takes the same value in all inertial frames of reference.

3.3.6 Transformation of Energy and Momentum

Consider two inertial reference frames, S and S' . Let S' move with velocity $\mathbf{v} = v \mathbf{e}_x$, and be in a standard configuration, with respect to S . Let \mathbf{p} and E be some particle's momentum and energy, respectively, in S . Likewise, let \mathbf{p}' and E' be the particle's momentum and energy, respectively, in S' . We have seen that the transformation of spacetime coordinates, (3.110)–(3.113), implies that the spacetime interval, $(c \Delta t)^2 - |\Delta \mathbf{r}|^2$, takes the same value in all inertial frames of reference. Given that $(E/c)^2 - |\mathbf{p}|^2$ also takes the same value in all inertial frames of reference, it seems reasonable to assume, by analogy, that the components of \mathbf{p} and E in our two inertial reference frames are related as follows:

$$p'_x = \gamma \left(p_x - \frac{v E}{c^2} \right), \quad (3.180)$$

$$p'_y = p_y, \quad (3.181)$$

$$p'_z = p_z, \quad (3.182)$$

$$E' = \gamma (E - v p_x). \quad (3.183)$$

We can easily test out the previous transformation rule. Suppose that the particle is at rest in S . It follows that $E = m_0 c^2$ and $p_x = p_y = p_z = 0$. Hence, Equations (3.180)–(3.183) yield

$$p'_x = -\gamma m_0 v = -m v, \quad (3.184)$$

$$p'_y = 0, \quad (3.185)$$

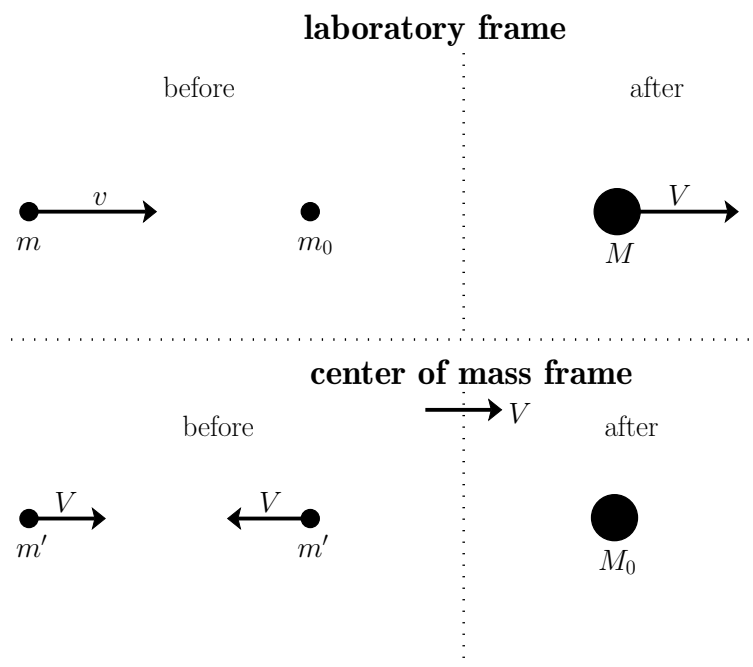


Figure 3.12: Relativistic inelastic collision.

$$p'_z = 0, \quad (3.186)$$

$$E' = \gamma m_0 c^2 = m c^2 \quad (3.187)$$

in S' . In other words, in the frame S' , in which the particle moves with velocity $-\mathbf{v}$, we have $\mathbf{p}' = -m \mathbf{v}$ and $E' = m c^2$. Of course, these are the correct results. (See Sections 3.3.2 and 3.3.4.)

3.3.7 Relativistic Momentum Conservation

Given that energy and momentum are clearly very closely related concepts in relativistic dynamics, we conjecture that conservation of energy also implies conservation of momentum.

Consider the situation illustrated in Figure 3.12. In the laboratory frame, a particle of rest mass m_0 , moving with speed v , collides with another particle of rest mass m_0 that is stationary. After the collision, the two particles stick together, and the composite particle, whose relativistic mass is M , moves off in the same direction as the originally moving particle at speed V . Now, the relativistic mass of the originally moving particle is

$$m = \frac{m_0}{\sqrt{1 - v^2/c^2}}. \quad (3.188)$$

Thus, if mass/energy is conserved in the collision then the net relativistic mass before the collision must match that after the collision. In other words,

$$m + m_0 = M, \quad (3.189)$$

which implies that

$$M = \left(1 + \frac{1}{\sqrt{1 - v^2/c^2}} \right) m_0. \quad (3.190)$$

On the other hand, if momentum is also conserved in the collision then the net relativistic momentum before the collision must match that after the collision. In other words,

$$m v = M V, \quad (3.191)$$

which yields

$$\frac{m_0}{\sqrt{1 - v^2/c^2}} = \left(1 + \frac{1}{\sqrt{1 - v^2/c^2}} \right) m_0 V, \quad (3.192)$$

where use has been made of Equations (3.188) and (3.190). The previous equation implies that

$$V = \frac{v}{1 + \sqrt{1 - v^2/c^2}}, \quad (3.193)$$

which can be rearranged to give

$$V = \frac{v - V}{1 - v V/c^2}. \quad (3.194)$$

However, there is another way of obtaining the previous equation. Let us transform to a frame of reference that moves, with respect to the laboratory frame, with speed V parallel to the motion of the original moving particle. The composite particle appears stationary in this reference frame. Now, if momentum is conserved, then the new reference frame is the center-of-mass frame. (See Section 1.6.1.) Consequently, our two particles must approach one another with equal and opposite velocities, V , before the collision, as shown in the figure. However, when the transformation of velocity, (3.122), is applied to the originally moving particle, we obtain

$$V = \frac{v - V}{1 - v V/c^2}, \quad (3.195)$$

which is identical to Equation (3.194). Note that we obtained Equation (3.194) from considerations of energy and momentum conservation, whereas we obtained the previous equation from a consideration of momentum conservation alone. Hence, we deduce that momentum conservation in relativistic dynamics implies energy conservation, and vice versa.

3.3.8 Photons

The general energy-momentum relation, (3.179), implies that a particle with zero rest mass has the simplified energy-momentum relation

$$E = p c. \quad (3.196)$$

Consider a photon. The *photoelectric effect* demonstrates that the energy of a photon is related to its angular frequency, ω , according to

$$E = \hbar \omega, \quad (3.197)$$

where \hbar is Planck's constant divided by 2π . (See Section 4.1.2.) However, we also know that a photon travels at the speed of light in all inertial reference frames. Thus, the relativistic energy, (3.177), of a photon can only be finite if the photon is a massless particle. In other words, a photon's rest mass must be zero. Hence, the previous two equations suggest that the momentum of a photon has the magnitude

$$p = \frac{\hbar \omega}{c}. \quad (3.198)$$

But, the dispersion relation of electromagnetic radiation in a vacuum, and, hence, of a photon moving through a vacuum, is

$$\omega = k c. \quad (3.199)$$

Here, \mathbf{k} is the wavevector of the radiation, and, hence, of the photon. Note that the direction of \mathbf{k} corresponds to the direction of motion of the photon. It is, thus, plausible that the momentum of our photon is written

$$\mathbf{p} = \hbar \mathbf{k}. \quad (3.200)$$

Consider the two inertial reference frames, S and S' , discussed in the Section 3.3.6. Let ω and \mathbf{k} be the angular frequency and wavevector, respectively, of our photon in S . Likewise, let ω' and \mathbf{k}' be the angular frequency and wavevector, respectively, of our photon in S' . Equations (3.180)–(3.183), (3.197), and (3.200) suggest that

$$k'_x = \gamma \left(k_x - \frac{v \omega}{c^2} \right), \quad (3.201)$$

$$k'_y = k_y, \quad (3.202)$$

$$k'_z = k_z, \quad (3.203)$$

$$\omega' = \gamma (\omega - v k_x). \quad (3.204)$$

3.3.9 Relativistic Doppler Effect

Consider the two inertial reference frames, S and S' , discussed in Section 3.3.6. Suppose that we place a radiation source at the origin of reference frame S . Let the source emit plane waves of angular frequency ω that travel in the positive x -direction. It follows that the wavevector of the radiation in S is $\mathbf{k} = (k_x, 0, 0)$, where

$$\omega = k_x c. \quad (3.205)$$

[See Equation (3.199).]

Consider an observer located at the origin of frame S' . To this observer, the radiation source appears to move at the speed $\mathbf{v} = -v \mathbf{e}_x$. Let $\mathbf{k}' = (k'_x, 0, 0)$ and ω' be the wavevector and angular frequency of the radiation measured by our observer. It follows from Equations (3.201), (3.204), and (3.205) that

$$k'_x = \gamma \left(k_x - \frac{v \omega}{c^2} \right) = \gamma \left(1 - \frac{v}{c} \right) k_x, \quad (3.206)$$

$$\omega' = \gamma (\omega - v k_x) = \gamma \left(1 - \frac{v}{c} \right) \omega. \quad (3.207)$$

Let $f = \omega/(2\pi)$ and $\lambda = 2\pi/k_x$ be the frequency (in hertz) and wavelength, respectively, of the radiation emitted by the source in its rest frame, and let $f' = \omega'/(2\pi)$ and $\lambda' = 2\pi/k'_x$ be the frequency (in hertz) and wavelength, respectively, of the radiation measured by an observer in the frame S' , in which the source recedes directly away from the observer at the speed v . Given that $\gamma = (1 - v^2/c^2)^{-1/2}$, we deduce that

$$\lambda' = \left(\frac{1 + v/c}{1 - v/c} \right)^{1/2} \lambda, \quad (3.208)$$

$$f' = \left(\frac{1 - v/c}{1 + v/c} \right)^{1/2} f. \quad (3.209)$$

It follows that if a radiation source recedes from an observer (or vice versa, because, in the absence of a medium through which electromagnetic waves propagate, all motion of sources and observers is relative) then the wavelength and frequency of the radiation measured by the observer will be larger and smaller, respectively, than the corresponding values measured in the rest frame of the source. This shift in the wavelength and frequency of electromagnetic radiation due to the relative motion of the observer and source is known as the *relativistic Doppler effect*. [Note that the non-relativistic Doppler effect for sound waves takes a different form to Equations (3.208) and (3.209) because motion of the source and the observer can be distinguished from one another in the presence of a medium through which the waves travel.]

By analogy with the previous two formulae, if the source moves directly toward the observer with the speed v then

$$\lambda' = \left(\frac{1 - v/c}{1 + v/c} \right)^{1/2} \lambda, \quad (3.210)$$

$$f' = \left(\frac{1 + v/c}{1 - v/c} \right)^{1/2} f. \quad (3.211)$$

In this case, the wavelength and frequency of the radiation measured by the observer are smaller and larger, respectively, than the corresponding values measured in the rest frame of the source. Thus, we can write the composite formulae

$$\lambda' = \left(\frac{1 \pm v/c}{1 \mp v/c} \right)^{1/2} \lambda, \quad (3.212)$$

$$f' = \left(\frac{1 \mp v/c}{1 \pm v/c} \right)^{1/2} f, \quad (3.213)$$

where the upper/lower signs correspond to the source moving directly away from/toward the observer (or vice versa).

3.3.10 Transverse Doppler Effect

Consider the situation illustrated in Figure 3.13. A radiation source that is located at the origin of reference frame S emits electromagnetic radiation of angular frequency ω , whose direction

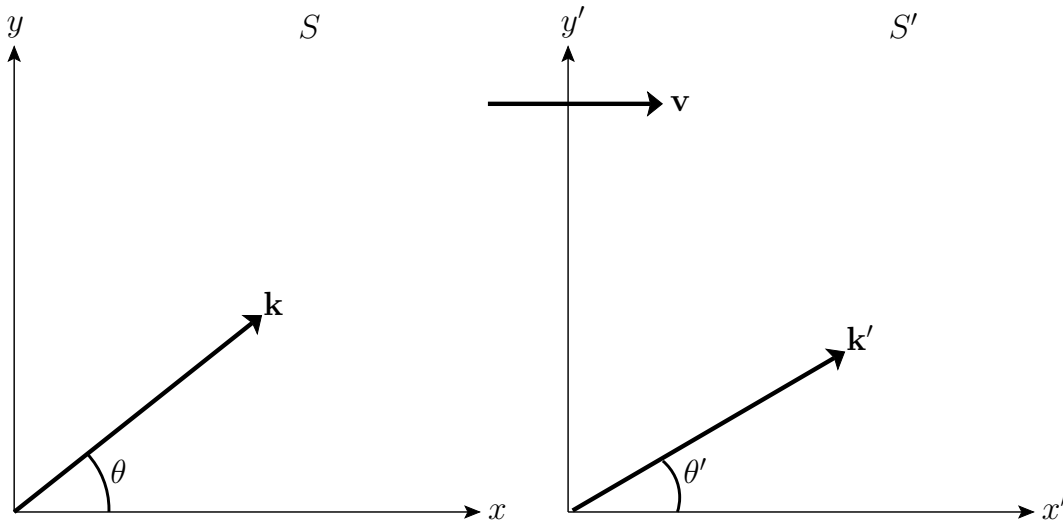


Figure 3.13: Relativistic Doppler effect.

of propagation lies in the x - y plane and subtends an angle θ with the x -axis. It follows that the wavevector of the radiation has the non-zero components

$$k_x = \frac{\omega}{c} \cos \theta, \quad (3.214)$$

$$k_y = \frac{\omega}{c} \sin \theta. \quad (3.215)$$

Suppose that the radiation is observed in a frame S' that moves with velocity $\mathbf{v} = v \mathbf{e}_x$, and is in a standard configuration, with respect to S . In S' , let ω' be the angular frequency of the radiation, and let θ' be the angle subtended by its direction of propagation and the x' -axis. It follows that, in S' , the wavevector of the radiation has the non-zero components

$$k'_x = \frac{\omega'}{c} \cos \theta', \quad (3.216)$$

$$k'_y = \frac{\omega'}{c} \sin \theta'. \quad (3.217)$$

The previous four equations can be combined with Equations (3.201)–(3.204) to give

$$\omega' \cos \theta' = \gamma \left(\cos \theta - \frac{v}{c} \right) \omega, \quad (3.218)$$

$$\omega' \sin \theta' = \omega \sin \theta, \quad (3.219)$$

$$\omega' = \gamma \left(1 - \frac{v}{c} \cos \theta \right) \omega. \quad (3.220)$$

Given that $\omega \propto f$, we deduce that

$$f' = \gamma \left(1 - \frac{v}{c} \cos \theta \right) f, \quad (3.221)$$

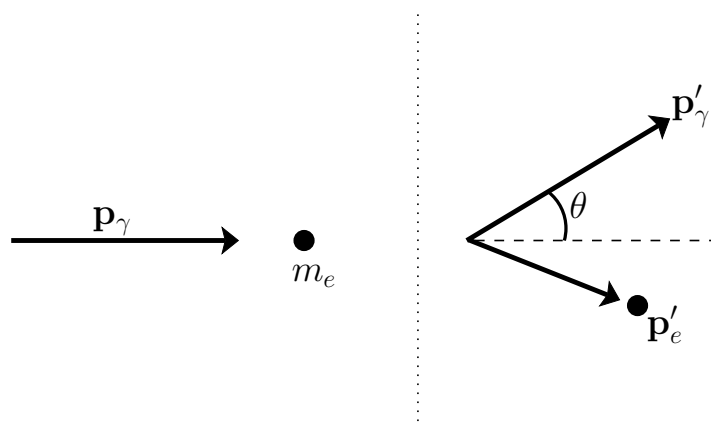


Figure 3.14: Compton Scattering.

which is the generalization of formula (3.209) to the case where the radiation is not propagating parallel to the relative velocity of the source and the observer. In particular, if $\theta = \pi/2$, so that the radiation is propagating in a direction that is perpendicular to the relative velocity of the source and the observer, then

$$f' = \gamma f. \quad (3.222)$$

In other words, in this case, the observer always measures an increased frequency of the radiation, relative to the frequency measured in the rest frame of the source. This effect is known as the *transverse Doppler effect*, and is a purely relativistic effect (i.e., it has no concomitant in Newtonian dynamics). The transverse Doppler effect was experimentally verified by Ives and Stilwell in 1938.

Finally, Equations (3.218) and (3.219) can be combined to give

$$\tan \theta' = \frac{\sin \theta}{\gamma (\cos \theta + v/c)}. \quad (3.223)$$

However, this formula is identical to the relativistic aberration formulae, (3.131), that we derived previously.

3.3.11 Compton Scattering

Compton scattering occurs when X-rays scatter off electrons in ordinary matter. The result is an increase in the wavelength of the scattered X-rays. This increase is inexplicable within the context of classical physics, which predicts that radiation that scatters off a stationary target should suffer no change in wavelength. In fact, as we shall explain, this effect can be explained in terms of the scattering of individual X-ray photons by individual electrons.

Consider the situation, illustrated in Figure 3.14, in which an X-ray photon of momentum \mathbf{p}_γ collides with a stationary electron of rest mass m_e . After the collision, the momentum of the photon is \mathbf{p}'_γ , and the recoil momentum of the electron is \mathbf{p}'_e . Conservation of momentum in the collision requires that

$$\mathbf{p}_\gamma = \mathbf{p}'_\gamma + \mathbf{p}'_e. \quad (3.224)$$

However, we know that

$$\mathbf{p}_\gamma = \hbar \mathbf{k}, \quad (3.225)$$

$$\mathbf{p}'_\gamma = \hbar \mathbf{k}', \quad (3.226)$$

$$\mathbf{p}'_e = \gamma m_e v, \quad (3.227)$$

where \mathbf{k} and \mathbf{k}' are the photon's initial and final wavevector, respectively, v is the electron's recoil speed, and $\gamma = (1 - v^2/c^2)^{-1/2}$. [See Equations (3.162) and (3.200).] Thus, we obtain

$$\mathbf{k} - \mathbf{k}' = \frac{m_e c}{\hbar} \gamma \frac{v}{c}. \quad (3.228)$$

The previous equation yields

$$|\mathbf{k} - \mathbf{k}'|^2 = \left(\frac{m_e c}{\hbar}\right)^2 \gamma^2 \frac{v^2}{c^2} = \left(\frac{m_e c}{\hbar}\right)^2 (\gamma^2 - 1), \quad (3.229)$$

or

$$k^2 - 2kk' \cos \theta + k'^2 = \left(\frac{m_e c}{\hbar}\right)^2 (\gamma^2 - 1). \quad (3.230)$$

Here, θ is the angle through which the photon is scattered (i.e., the angle subtended between \mathbf{k} and \mathbf{k}'). See Figure 3.14.

Let E_γ , E'_γ , E_e , and E'_e be the initial photon energy, the final photon energy, the initial electron energy, and the final electron energy, respectively. Energy conservation in the collision requires that

$$E_\gamma + E_e = E'_\gamma + E'_e. \quad (3.231)$$

However, we know that

$$E_\gamma = \hbar c k, \quad (3.232)$$

$$E'_\gamma = \hbar c k', \quad (3.233)$$

$$E_e = m_e c^2, \quad (3.234)$$

$$E'_e = \gamma m_e c^2. \quad (3.235)$$

[See Equations (3.174), (3.197), and (3.199).] Hence, we get

$$\gamma = \frac{k - k' + m_e c / \hbar}{m_e c / \hbar}. \quad (3.236)$$

Equations (3.230) and (3.236) can be combined to give

$$k^2 - 2kk' \cos \theta + k'^2 = \left(k - k' + \frac{m_e c}{\hbar}\right)^2 - \left(\frac{m_e c}{\hbar}\right)^2, \quad (3.237)$$

or

$$k^2 - 2kk' \cos \theta + k'^2 = k^2 - 2kk' + k'^2 + 2(k - k') \frac{m_e c}{\hbar}, \quad (3.238)$$

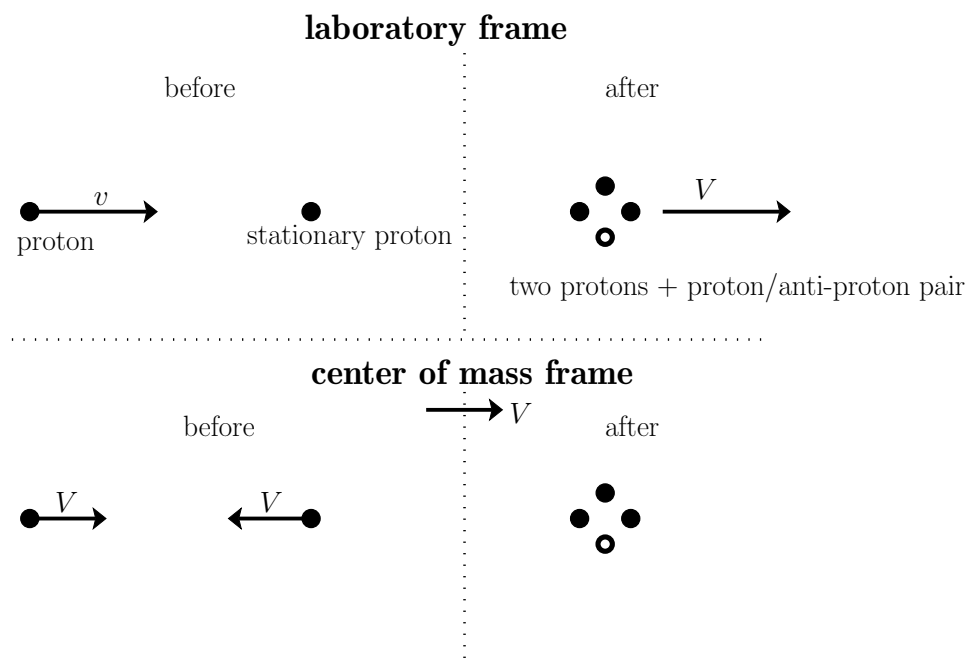


Figure 3.15: Pair creation.

which can be rearranged to produce

$$\frac{1}{k'} - \frac{1}{k} = \frac{\hbar}{m_e c} (1 - \cos \theta). \quad (3.239)$$

Finally, if $\lambda = 2\pi/k$ and $\lambda' = 2\pi/k'$ are the initial and final wavelengths of the photon then we obtain

$$\lambda' - \lambda = \frac{h}{m_e c} (1 - \cos \theta). \quad (3.240)$$

The previous equation relates the increase in wavelength of the scattered photon to its scattering angle in a simple manner. Here, $h/(m_e c) = 2.43 \times 10^{-12}$ m is known as the *Compton wavelength* of the electron. The previous formula was verified experimentally by Arthur Compton in 1923.

3.3.12 Relativistic Inelastic Scattering

Finally, consider the situation, illustrated in Figure 3.15, in which a moving proton collides with a stationary proton, and a proton/anti-proton pair is created during the collision. We wish to determine the minimum energy of the incident proton required to create the pair. Let m_p be the proton rest mass. As is clear from the figure, in the center of mass frame, the minimum energy state corresponds to the case in which the particles are all at rest after the collision. (Additional energy would just causes the particles to move away from one another, in this frame, after the collision.) Thus, in the laboratory frame, the particles must all move with a common velocity after the collision. However, given that the particles all have the same mass, each particle in the laboratory frame

must have momentum $P/4$, after the collision, where P is the total momentum of the system, and must have energy $E/4$, where E is the total energy of the system. Thus, the energy-momentum relation [see Equation (3.179)] for one of the particles after the collision yields

$$\left(\frac{E}{4}\right)^2 = m_p^2 c^4 + \left(\frac{Pc}{4}\right)^2, \quad (3.241)$$

or

$$E^2 = 16 m_p^2 c^4 + P^2 c^2. \quad (3.242)$$

Let E_0 be the initial laboratory-frame energy of the moving proton before the collision. Energy conservation requires that

$$E_0 + m_p c^2 = E. \quad (3.243)$$

The previous two equations can be combined to give

$$(E_0 + m_p c^2)^2 = 16 m_p^2 c^4 + P^2 c^2, \quad (3.244)$$

or

$$E_0^2 + 2 E_0 m_p c^2 = 15 m_p^2 c^4 + P^2 c^2. \quad (3.245)$$

However, the initial momentum of the moving proton in the laboratory frame is P (because the proton possesses all of the system's initial momentum, and the total momentum must be the same before and after the collision). Hence, the moving proton's initial energy-momentum relation [see Equation (3.179)] is

$$E_0^2 = m_p^2 c^4 + P^2 c^2. \quad (3.246)$$

The previous two equations yield

$$2 E_0 m_p c^2 = 14 m_p^2 c^4, \quad (3.247)$$

or

$$E_0 = 7 m_p c^2. \quad (3.248)$$

Thus, the minimum kinetic energy of the incident proton required to generate a proton/anti-proton pair is 6 times its rest mass energy. This corresponds to a Lorentz factor of 7, which implies a speed of about 99% of the speed of light in vacuum.

3.4 Relativity and Electromagnetism

3.4.1 Transformation of Electromagnetic Fields

In this section, we shall investigate how electromagnetic fields transform when viewed in different inertial frames of reference. Our investigation is premised on two assumptions. First, Maxwell's equations (see Section 2.4.2) take equivalent forms in all inertial frames of reference. Of course, this is just a special case of the equivalence principle discussed in Section 3.2.1. Second, the

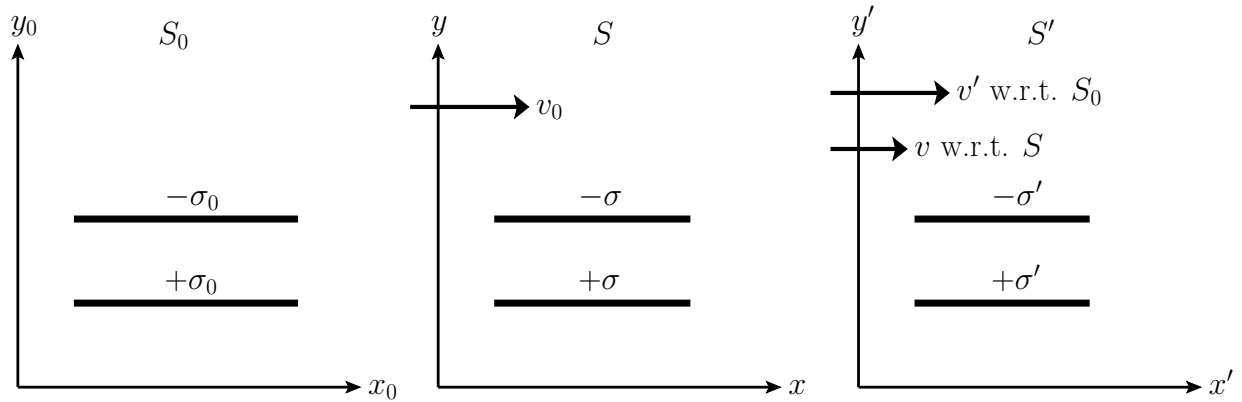


Figure 3.16: Parallel plate capacitor.

electric charge of an elementary particle is the same in all inertial reference frames. Our second assumption is an experimentally verifiable fact.

Consider three inertial reference frames, S_0 , S , and S' , that are all in standard configurations with respect to one another. (See Section 3.2.6.) Let frame S move parallel to the x -axis at speed v_0 with respect to frame S_0 . Let frame S' move parallel to the x -axis at speed v' with respect to frame S_0 , and with speed v with respect to frame S . See Figure 3.16. It follows from the relativistic transformation of velocity (see Section 3.2.9) that

$$v' = \frac{v_0 + v}{1 + v_0 v/c^2}. \quad (3.249)$$

Let us define the Lorentz factors

$$\gamma_0 = \left(1 - \frac{v_0^2}{c^2}\right)^{-1/2}, \quad (3.250)$$

$$\gamma = \left(1 - \frac{v^2}{c^2}\right)^{-1/2}, \quad (3.251)$$

$$\gamma' = \left(1 - \frac{v'^2}{c^2}\right)^{-1/2}. \quad (3.252)$$

The previous four equations yield

$$\begin{aligned} \gamma' &= \left[1 - \frac{1}{c^2} \left(\frac{v_0 + v}{1 + v_0 v/c^2}\right)^2\right]^{-1/2} \\ &= \left(1 + \frac{v_0 v}{c^2}\right) \left[\left(1 + \frac{v_0 v}{c^2}\right)^2 - \left(\frac{v_0}{c} + \frac{v}{c}\right)^2\right]^{-1/2} \\ &= \left(1 + \frac{v_0 v}{c^2}\right) \left(1 - \frac{v_0^2}{c^2} - \frac{v^2}{c^2} + \frac{v_0^2 v^2}{c^4}\right)^{-1/2} \end{aligned}$$

$$\begin{aligned}
&= \left(1 + \frac{v_0 v}{c^2}\right) \left(1 - \frac{v_0^2}{c^2}\right)^{-1/2} \left(1 - \frac{v^2}{c^2}\right)^{-1/2} \\
&= \gamma_0 \gamma \left(1 + \frac{v_0 v}{c^2}\right). \tag{3.253}
\end{aligned}$$

Suppose that frame S_0 contains a parallel plate capacitor that is at rest. Let the capacitor plates be parallel to the x - z plane. Furthermore, let the lower (in y) plate have the uniform electric charge density σ_0 , and let the upper plate have the uniform charge density $-\sigma_0$. See Figure 3.16. In frame S , the capacitor plates appear to move in the $-x$ -direction with speed v_0 . Thus, the lengths of the plates (in the x -direction) are contracted by a factor γ_0 , whereas the widths of the plates (in the z -direction) are unchanged. (See Section 3.2.7.) Moreover, according to our second assumption, the net electric charges on the two capacitor plates in frame S are the same as those in frame S_0 . It follows that, in frame S , the charge densities on the two plates are $\pm\sigma$, where

$$\sigma = \gamma_0 \sigma_0. \tag{3.254}$$

Analogous reasoning reveals that the charge densities on the two capacitor plates in frame S' are $\pm\sigma'$, where

$$\sigma' = \gamma' \sigma_0. \tag{3.255}$$

All of the electric charges in frame S_0 are stationary, so the associated current density is zero. However, in frames S and S' , the charges on the capacitor plates appear to move in the $-x$ -direction with speeds v_0 and v' , respectively. Thus, in frame S , the current per unit width (in the z -direction) flowing on the lower (in y) capacitor plate takes the form $\mathbf{J} = J_x \mathbf{e}_x$, where

$$J_x = -\sigma v_0 = -\gamma_0 v_0 \sigma_0. \tag{3.256}$$

There is an equal and opposite current per unit width flowing on the upper plate. Likewise, in frame S' , the current per unit width flowing on the lower plate takes the form $\mathbf{J}' = J'_x \mathbf{e}_x$, where

$$J'_x = -\sigma' v' = -\gamma' v' \sigma_0. \tag{3.257}$$

Again, there is an equal and opposite current per unit width flowing on the upper plate.

The integral form of the Maxwell equation (2.484) is

$$\oint_S \mathbf{E} \cdot d\mathbf{S} = \frac{1}{\epsilon_0} \int_V \rho dV, \tag{3.258}$$

where S is some surface enclosing a volume V , and where use has been made of the divergence theorem. (See Section A.20.) As described in Sections 2.1.12 and 2.1.13, if the previous equation is applied to a Gaussian pill-box in frame S that encloses one or other of the capacitor plates, and co-moves with the plates, then it is easily demonstrated that the electric field in the region between the plates is uniform, taking the form $\mathbf{E} = E_y \mathbf{e}_y$, where

$$E_y = \frac{\sigma}{\epsilon_0} = \gamma_0 \frac{\sigma_0}{\epsilon_0}. \tag{3.259}$$

Note, incidentally, that there is nothing in Equation (3.258) that precludes volume V from being a moving volume. Analogous reasoning reveals that the electric field between the capacitor plates in frame S' is uniform, taking the value $\mathbf{E}' = E'_y \mathbf{e}_y$, where

$$E'_y = \frac{\sigma'}{\epsilon_0} = \gamma' \frac{\sigma_0}{\epsilon_0}. \quad (3.260)$$

The integral form of the Maxwell equation (2.487) is

$$\oint_C \mathbf{B} \cdot d\mathbf{r} = \int_S \left(\mu_0 \mathbf{j} + \epsilon_0 \mu_0 \frac{\partial \mathbf{E}}{\partial t} \right) \cdot d\mathbf{S}, \quad (3.261)$$

where S is a surface bounded by a loop C , and where use has been made of the curl theorem. (See Section A.22.) However, the electric field between the capacitor plates is constant in time in all three of our reference frames, so the previous equation simplifies to give

$$\oint_C \mathbf{B} \cdot d\mathbf{r} = \mu_0 \int_S \mathbf{j} \cdot d\mathbf{S}, \quad (3.262)$$

If the previous equation is applied to an Ampèrian loop in the y - z plane of frame S that straddles one or other of the capacitor plates, and co-moves with the plates, then it is easily demonstrated that the magnetic field in the region between the plates is uniform, taking the form $\mathbf{B} = B_z \mathbf{e}_z$, where

$$B_z = \mu_0 J_x = -\gamma_0 v_0 \mu_0 \sigma_0. \quad (3.263)$$

Again, there is nothing in Equation (3.262) that prohibits surface S from being a moving surface. Likewise, in frame S' , the magnetic field between the plates is uniform, taking the form $\mathbf{B} = B'_z \mathbf{e}_z$, where

$$B'_z = \mu_0 J'_x = -\gamma' v' \mu_0 \sigma_0. \quad (3.264)$$

According to Equations (3.253), (3.259), (3.260), and (3.263),

$$\begin{aligned} E'_y &= \gamma_0 \gamma \left(1 + \frac{v_0 v}{c^2} \right) \frac{\sigma_0}{\epsilon_0} \\ &= \gamma \left(\gamma_0 \frac{\sigma_0}{\epsilon_0} + v \gamma_0 v_0 \mu_0 \sigma_0 \right) \\ &= \gamma (E_y - v B_z), \end{aligned} \quad (3.265)$$

where use has been made of $c = 1/\sqrt{\epsilon_0 \mu_0}$. Likewise, Equation (3.249), (3.253), (3.259), (3.263), and (3.264) yield

$$\begin{aligned} B'_z &= -\gamma_0 \gamma \left(1 + \frac{v_0 v}{c^2} \right) \left(\frac{v_0 + v}{1 + v_0 v/c^2} \right) \mu_0 \sigma_0 \\ &= \gamma \left(-\gamma_0 v_0 \mu_0 \sigma_0 - \frac{v}{c^2} \gamma_0 \frac{\sigma_0}{\epsilon_0} \right) \end{aligned}$$

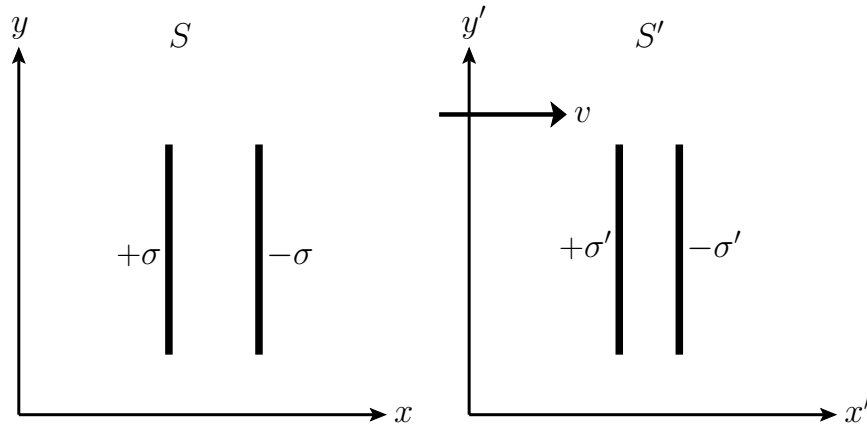


Figure 3.17: Parallel plate capacitor.

$$= \gamma \left(B_z - \frac{v}{c^2} E_y \right). \quad (3.266)$$

If we repeat the previous exercise with capacitor plates that are parallel to the x - y plane, instead of the x - z plane, then it is easily demonstrated that

$$E'_z = \gamma (E_z + v B_y), \quad (3.267)$$

$$B'_y = \gamma \left(B_y + \frac{v}{c^2} E_z \right). \quad (3.268)$$

In order to determine the transformation rule for E_x , consider the situation shown in Figure 3.17. Here, a parallel plate capacitor is stationary in frame S , and is aligned such that its plates are parallel to the y - z plane. Frame S' is in a standard configuration, and moves with velocity $\mathbf{v} = v \mathbf{e}_x$, with respect to frame S . The cross-sectional areas of the plates are the same in reference frames S and S' . (See Section 3.2.7.) Furthermore, the electric charges on the plates are identical in both frames. Hence, we deduce that the electric charge densities on the plates are the same in the two reference frames. In other words,

$$\sigma' = \sigma. \quad (3.269)$$

Now, the parallel distance between the two plates is contracted by a factor γ in frame S' , compared to frame S . However, the electric field generated between the capacitor plates only depends on the charge density residing on the plates, and is independent of the inter-plate spacing. In fact, the electric fields in frames S and S' are $\mathbf{E} = E_x \mathbf{e}_x$ and $\mathbf{E}' = E'_x \mathbf{e}_x$, respectively, where

$$E_x = \frac{\sigma}{\epsilon_0}, \quad (3.270)$$

$$E'_x = \frac{\sigma'}{\epsilon_0}. \quad (3.271)$$

Hence, we deduce from the previous three equations that

$$E'_x = E_x. \quad (3.272)$$

In order to determine the transformation for B_x , consider a long, thin, solenoid whose axis runs parallel to the x -direction. Let the solenoid have N turns per unit length, and carry a current I , in frame S , and let the solenoid have N' turns per unit length, and carry a current I' , in frame S' . Frame S' is in a standard configuration, and moves with velocity $\mathbf{v} = v \mathbf{e}_x$, with respect to frame S . Because of relativistic length contraction (see Section 3.2.4),

$$N' = \gamma N. \quad (3.273)$$

On the other hand, because current is electric charge per unit time, and electric charge is invariant between different inertial frames, whereas time is dilated in frame S' , relative to frame S , we have

$$I' = \frac{I}{\gamma}. \quad (3.274)$$

According to Section 2.2.11, the magnetic fields generated inside the solenoid in frames S and S' are $\mathbf{B} = B_x \mathbf{e}_x$ and $\mathbf{B}' = B'_x \mathbf{e}_x$, respectively, where

$$B_x = \mu_0 N I, \quad (3.275)$$

$$B'_x = \mu_0 N' I'. \quad (3.276)$$

The previous four equations imply that

$$B'_x = B_x. \quad (3.277)$$

Thus, we can now state the complete set of transformation laws for the components of the electric and magnetic field between an inertial reference frame S , and a second inertial reference frame, S' , that is in a standard configuration, and moves at velocity $\mathbf{v} = v \mathbf{e}_x$, with respect to the first. The transformation laws are as follows:

$$E'_x = E_x, \quad (3.278)$$

$$E'_y = \gamma(E_y - v B_z), \quad (3.279)$$

$$E'_z = \gamma(E_z + v B_y), \quad (3.280)$$

and

$$B'_x = B_x, \quad (3.281)$$

$$B'_y = \gamma \left(B_y + \frac{v}{c^2} E_z \right), \quad (3.282)$$

$$B'_z = \gamma \left(B_z - \frac{v}{c^2} E_y \right). \quad (3.283)$$

Here, $\gamma = (1 - v^2/c^2)^{-1/2}$.

It is easily demonstrated from the transformation rules (3.278)–(3.283) that

$$\mathbf{E}' \cdot \mathbf{B}' = \mathbf{E} \cdot \mathbf{B}, \quad (3.284)$$

$$E'^2 - c^2 B'^2 = E^2 - c^2 B^2. \quad (3.285)$$

Thus, if electric and magnetic fields in one inertial frame of reference are in the configuration of an electromagnetic wave traveling through a vacuum—in other words, if $\mathbf{E} \cdot \mathbf{B} = 0$ and $E = c B$ (see Section 2.4.4)—then they are in this configuration in all inertial frames of reference.

3.4.2 Electromagnetic Fields of a Moving Charge

Consider an electric charge, q , that is at rest at the origin of an inertial reference frame S' . The electric and magnetic fields generated by the charge at displacement \mathbf{r}' in frame S' are

$$\mathbf{E}' = \frac{q}{4\pi\epsilon_0} \frac{\mathbf{r}'}{r'^3}, \quad (3.286)$$

$$\mathbf{B}' = \mathbf{0}, \quad (3.287)$$

respectively. (See Sections 2.1.2 and 2.2.7.) Thus,

$$E'_x = \frac{q}{4\pi\epsilon_0} \frac{x'}{(x'^2 + y'^2 + z'^2)^{3/2}}, \quad (3.288)$$

$$E'_y = \frac{q}{4\pi\epsilon_0} \frac{y'}{(x'^2 + y'^2 + z'^2)^{3/2}}, \quad (3.289)$$

$$E'_z = \frac{q}{4\pi\epsilon_0} \frac{z'}{(x'^2 + y'^2 + z'^2)^{3/2}}, \quad (3.290)$$

$$B'_x = B'_y = B'_z = 0. \quad (3.291)$$

Let us transform to an inertial reference frame, S , that is in a standard configuration, and moves with velocity $-v\mathbf{e}_x$, with respect to frame S' . Thus, in frame S , the charge appears to move with velocity $\mathbf{v} = v\mathbf{e}_x$. Making use of the field transformation relations, (3.278)–(3.283), with primed and unprimed fields swapped, and $v \rightarrow -v$, we obtain

$$E_x = E'_x = \frac{q}{4\pi\epsilon_0} \frac{x'}{(x'^2 + y'^2 + z'^2)^{3/2}}, \quad (3.292)$$

$$E_y = \gamma(E'_y + vB'_z) = \frac{q}{4\pi\epsilon_0} \frac{\gamma y'}{(x'^2 + y'^2 + z'^2)^{3/2}}, \quad (3.293)$$

$$E_z = \gamma(E'_z - vB'_y) = \frac{q}{4\pi\epsilon_0} \frac{\gamma z'}{(x'^2 + y'^2 + z'^2)^{3/2}}, \quad (3.294)$$

$$B_x = B'_x = 0, \quad (3.295)$$

$$B_y = \gamma\left(B'_y - \frac{v}{c^2}E'_z\right) = \frac{q}{4\pi\epsilon_0} \frac{\gamma z'}{(x'^2 + y'^2 + z'^2)^{3/2}} \left(-\frac{v}{c^2}\right), \quad (3.296)$$

$$B_z = \gamma\left(B'_z + \frac{v}{c^2}E'_y\right) = \frac{q}{4\pi\epsilon_0} \frac{\gamma y'}{(x'^2 + y'^2 + z'^2)^{3/2}} \left(+\frac{v}{c^2}\right). \quad (3.297)$$

Consider the electric and magnetic fields generated by the charge at some point P in frame S whose displacement is $\mathbf{r} = (x, y, z)$. See Figure 3.18. The displacement of the charge in frame S is $\mathbf{r}' = (vt, 0, 0)$. Let

$$\mathbf{s} = \mathbf{r} - \mathbf{r}' = (x - vt, y, z) \quad (3.298)$$

be a vector that is directed from the instantaneous position of the charge in frame S to point P . A Lorentz transformation (see Section 3.2.7) between frames S and S' reveals that

$$x' = \gamma(x - vt) = \gamma s_x, \quad (3.299)$$

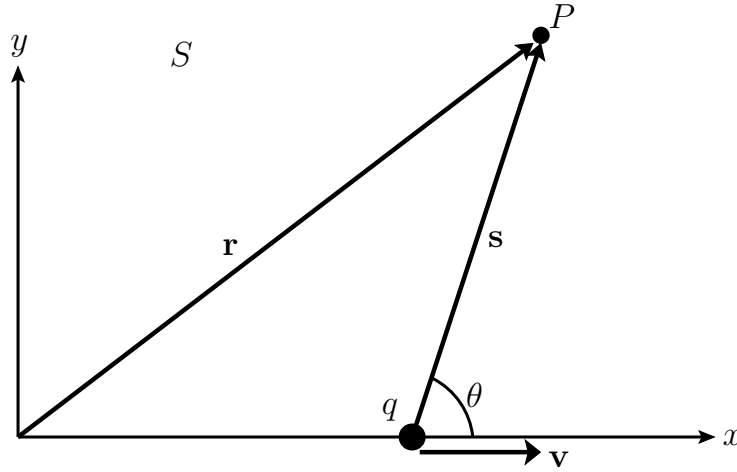


Figure 3.18: Observing the field of a moving charge.

$$y' = y = s_y, \quad (3.300)$$

$$z' = z = s_z. \quad (3.301)$$

Let us write

$$s_x = s \cos \theta, \quad (3.302)$$

$$s_y = s \sin \theta \cos \varphi, \quad (3.303)$$

$$s_z = s \sin \theta \sin \varphi, \quad (3.304)$$

where θ is the angle subtended between \mathbf{s} and \mathbf{v} , and φ is an azimuthal angle. See Figure 3.18. It is easily demonstrated from the previous six equations that

$$\begin{aligned} x'^2 + y'^2 + z'^2 &= \gamma^2 s_x^2 + s_y^2 + s_z^2 \\ &= (\gamma^2 \cos^2 \theta + \sin^2 \theta) s^2 = [(\gamma^2 - 1) \cos^2 \theta + 1] s^2 \\ &= \left(\frac{\gamma^2 v_r^2}{c^2} + 1 \right) s^2, \end{aligned} \quad (3.305)$$

where $v_r = v \cos \theta$ is the component of \mathbf{v} that is directed from the instantaneous position of the charge to the point P . Thus, making use of Equations (3.292)–(3.297) and Equations (3.299)–(3.301), we obtain

$$E_x = \frac{q}{4\pi \epsilon_0} \frac{\gamma}{(1 + \gamma^2 v_r^2/c^2)^{3/2}} \frac{s_x}{s^3}, \quad (3.306)$$

$$E_y = \frac{q}{4\pi \epsilon_0} \frac{\gamma}{(1 + \gamma^2 v_r^2/c^2)^{3/2}} \frac{s_y}{s^3}, \quad (3.307)$$

$$E_z = \frac{q}{4\pi \epsilon_0} \frac{\gamma}{(1 + \gamma^2 v_r^2/c^2)^{3/2}} \frac{s_z}{s^3}, \quad (3.308)$$

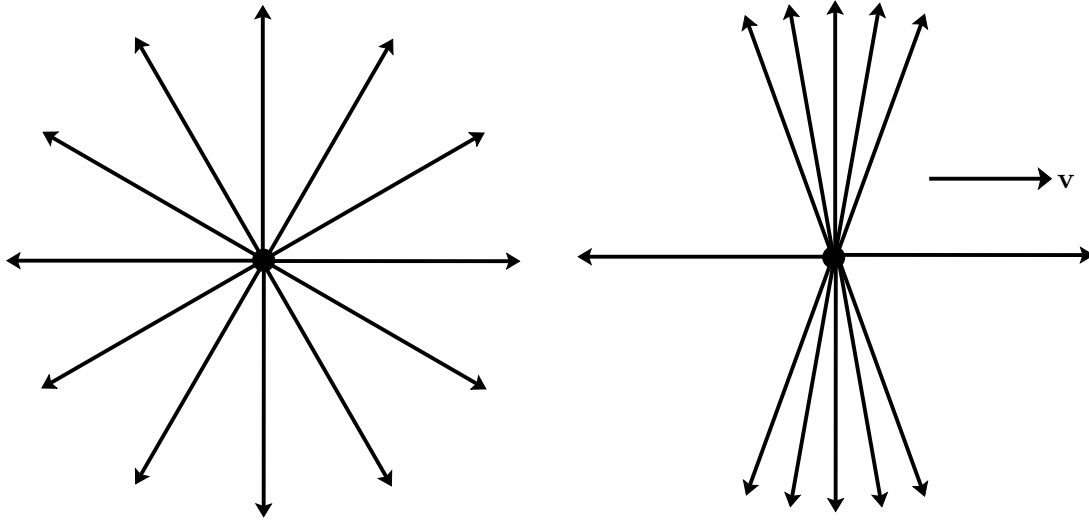


Figure 3.19: Electric field-lines of a stationary (left) and a moving (right) electric charge.

$$B_x = 0, \quad (3.309)$$

$$B_y = \frac{q}{4\pi\epsilon_0} \frac{\gamma}{(1 + \gamma^2 v_r^2/c^2)^{3/2}} \frac{s_z}{s^3} \left(-\frac{v}{c^2} \right), \quad (3.310)$$

$$B_z = \frac{q}{4\pi\epsilon_0} \frac{\gamma}{(1 + \gamma^2 v_r^2/c^2)^{3/2}} \frac{s_y}{s^3} \left(+\frac{v}{c^2} \right). \quad (3.311)$$

Note that the electric field-lines generated by a moving electric charge are straight-lines that are directed from the instantaneous position of the charge to the point of observation. At low velocities (i.e., $v \ll c$), the field-lines are equally spaced around the charge. However, as $v \rightarrow c$, the field-lines become increasingly bunched in the plane transverse to the charge's direction of motion that passes through the charge. This is illustrated schematically in Figure 3.19. The magnetic field generated by a moving charge is

$$\mathbf{B} = \frac{\mathbf{v} \times \mathbf{E}}{c^2}, \quad (3.312)$$

where \mathbf{v} is the charge's velocity, and \mathbf{E} is the electric field generated by the charge.

Chapter 4

Quantum Mechanics

4.1 Experimental Basis of Quantum Mechanics

4.1.1 Wave-Particle Duality

According to classical physics (i.e., physics prior to the 20th century), particles and waves are distinct classes of physical entities that possess markedly different properties. For instance, particles are discrete, which implies that they cannot be arbitrarily divided. In other words, it makes sense to talk about one electron, or two electrons, but not about a third of an electron. Waves, on the other hand, are continuous, which implies that they can be arbitrarily divided. In other words, given a wave whose amplitude has a certain value, it makes sense to talk about a similar wave whose amplitude is one third, or any other fraction whatsoever, of this value. Particles are also highly localized in space. For example, atomic nuclei have very small radii of order 10^{-15} m, whereas electrons act like point particles (i.e., they have no discernible spatial extent). Waves, on the other hand, are non-localized in space. In fact, a wave is defined as a disturbance that is periodic in space, with some finite periodicity length (i.e., wavelength). Hence, it is fairly meaningless to talk about a disturbance being a wave unless it extends over a region of space that is at least a few wavelengths in size.

The classical scenario, just described, in which particles and waves are distinct phenomena, had to be significantly modified in the early decades of the 20th century. During this time period, physicists discovered, much to their surprise, that, under certain circumstances, waves act as particles, and particles act as waves. This bizarre behavior is known as *wave-particle duality*. For instance, the *photoelectric effect* (see Section 4.1.2) shows that electromagnetic waves sometimes act like swarms of massless particles called *photons*. Moreover, the phenomenon of *electron diffraction* by atomic lattices (see Section 4.1.6) implies that electrons sometimes possess wave-like properties.

Wave-particle duality usually only manifests itself on atomic and sub-atomic lengthscales (i.e., on lengthscales less than, or of order, 10^{-10} m; see Section 4.1.6.) The classical picture remains valid on significantly longer lengthscales. Thus, on macroscopic lengthscales, waves only act like waves, particles only act like particles, and there is no wave-particle duality. However, on atomic lengthscales, classical mechanics, which governs the macroscopic behavior of massive particles, and classical electrodynamics, which governs the macroscopic behavior of electromagnetic

fields—neither of which take wave-particle duality into account—must be replaced by new theories. The theories in question are called *quantum mechanics* and *quantum electrodynamics*, respectively. In this section, we shall discuss a simple version of quantum mechanics in which the microscopic dynamics of massive particles (i.e., particles with finite mass) is described entirely in terms of wavefunctions. This particular version of quantum mechanics is known as *wave mechanics*. But, first, let us discuss the experimental evidence for wave-particle duality in more detail.

4.1.2 Photoelectric Effect

The so-called *photoelectric effect*, by which a polished metal surface emits electrons when illuminated by visible or ultra-violet light, was discovered by Heinrich Hertz in 1887. The following facts regarding this effect can be established via careful observation. First, a given surface only emits electrons when the frequency of the light with which it is illuminated exceeds a certain threshold value that is a property of the metal. Second, the current of photoelectrons, when it exists, is proportional to the intensity of the light falling on the surface. Third, the energy of the photoelectrons is independent of the light intensity, but varies linearly with the light frequency. These facts are inexplicable within the framework of classical physics.

In 1905, Albert Einstein proposed a radical new theory of light in order to account for the photoelectric effect. According to this theory, light of fixed angular frequency ω consists of a collection of indivisible discrete packages, called *quanta*,¹ whose energy is

$$E = \hbar \omega. \quad (4.1)$$

Here,

$$\hbar = 1.055 \times 10^{-34} \text{ J s} \quad (4.2)$$

is a new constant of nature, known as *Planck's constant*. (Strictly speaking, it is Planck's constant, h , divided by 2π .) Incidentally, \hbar is called Planck's constant, rather than Einstein's constant, because Max Planck first introduced the concept of the quantization of light, in 1900, while trying to account for the electromagnetic spectrum of a black body (i.e., a perfect emitter and absorber of electromagnetic radiation). (See Section 5.6.2.)

Suppose that the electrons at the surface of a piece of metal lie in a potential well of depth W . In other words, the electrons have to acquire an energy W in order to be emitted from the surface. Here, W is generally called the *workfunction* of the surface, and is a property of the metal. Suppose that an electron absorbs a single quantum of light, otherwise known as a *photon*. Its energy therefore increases by $\hbar \omega$. If $\hbar \omega$ is greater than W then the electron is emitted from the surface with the residual kinetic energy

$$K = \hbar \omega - W. \quad (4.3)$$

Otherwise, the electron remains trapped in the potential well, and is not emitted. Here, we are assuming that the probability of an electron absorbing two or more photons is negligibly small compared to the probability of it absorbing a single photon (as is, indeed, the case for relatively

¹Plural of *quantum*: Latin neuter of *quantus*: how much?

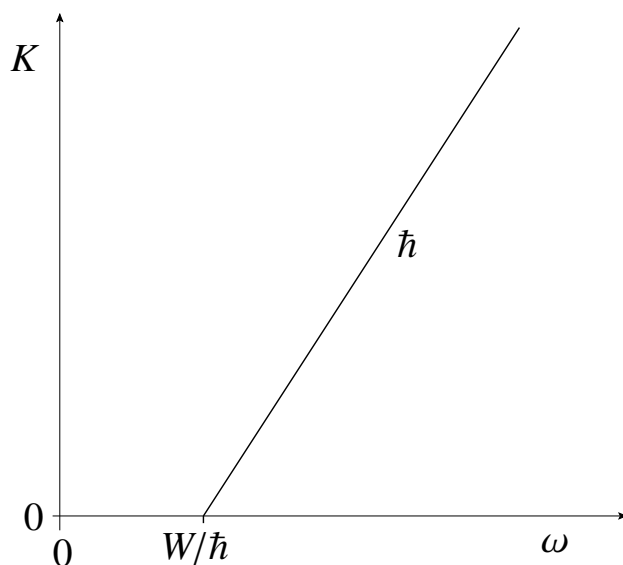


Figure 4.1: Variation of the kinetic energy, K , of photoelectrons with the wave angular frequency, ω .

low intensity illumination). Incidentally, we can determine Planck's constant, as well as the work-function of the metal, by plotting the kinetic energy of the emitted photoelectrons as a function of the wave frequency, as shown in Figure 4.1. This plot is a straight line whose slope is \hbar , and whose intercept with the ω axis is W/\hbar . Finally, the number of emitted electrons increases with the intensity of the light because, the more intense the light, the larger the flux of photons onto the surface. Thus, Einstein's quantum theory of light is capable of accounting for all three of the previously mentioned observational facts regarding the photoelectric effect.

Of course, an electromagnetic wave of angular frequency ω propagates through a vacuum at the speed of light in vacuum, c . However, if such a wave actually consists of a swarm of photons then it seems reasonable to suppose that these photons also move through a vacuum at the speed c . As discussed in Section 3.3.8, if photons move at the speed c then Einstein's special theory of relativity demands that they be *massless* particles with momenta

$$\mathbf{p} = \hbar \mathbf{k}, \quad (4.4)$$

where \mathbf{k} is the wavenumber of the associated electromagnetic wave.

4.1.3 Compton Scattering

As described in Section 3.3.11, formulae (4.1) and (4.4) for the energy and momentum, respectively, of a photon were directly verified experimentally when the phenomenon of *Compton scattering* was discovered in 1923.

4.1.4 Photon Polarization

In 1924, Frank Bubb discovered that if plane-polarized light is used to eject photo-electrons then there is a preferred direction of emission of the electrons. Clearly, the polarization properties of light, which are usually associated with its wave-like behavior, also extend to its particle-like behavior. In particular, a polarization can be ascribed to each individual photon in a beam of light.

Consider the following well-known experiment. A beam of plane polarized light is passed through a thin polarizing film whose plane is normal to the beam's direction of propagation, and which has the property that it is only transparent to light whose direction of polarization lies perpendicular to its optic axis (which is assumed to lie in the plane of the film). Classical electromagnetic wave theory tells us that if the beam is polarized perpendicular to the optic axis then all of the light is transmitted, if the beam is polarized parallel to the optic axis then none of the light is transmitted, and if the light is polarized at an angle α to the axis then a fraction $\sin^2 \alpha$ of the beam energy is transmitted; the latter result is known as *Malus's law*, after Étienne-Louis Malus who discovered it in 1808. Let us try to account for these observations at the individual photon level.

A beam of light that is plane polarized in a certain direction is presumably made up of a stream of photons that are each plane polarized in that direction. This picture leads to no difficulty if the direction of polarization lies parallel or perpendicular to the optic axis of the polarizing film. In the former case, none of the photons are transmitted, and, in the latter case, all of the photons are transmitted. But, what happens in the case of an obliquely polarized incident beam?

The previous question is not very precise. Let us reformulate it as a question relating to the result of some experiment that we could perform. Suppose that we were to fire a single photon at a polarizing film, and then look to see whether or not it emerges on the other side. The possible results of the experiment are that either a whole photon (whose energy is equal to the energy of the incident photon) is observed, or no photon is observed. Any photon that is transmitted through the film must be polarized perpendicular to the film's optic axis. Furthermore, it is impossible to imagine (in physics) finding part of a photon on the other side of the film. If we repeat the experiment a great number of times then, on average, a fraction $\sin^2 \alpha$ of the photons are transmitted through the film, and a fraction $\cos^2 \alpha$ are absorbed. Thus, given that the trials are statistically independent of one another, we must conclude that an individual photon has a probability $\sin^2 \alpha$ of being transmitted as a photon polarized in the plane perpendicular to the optic axis, and a probability $\cos^2 \alpha$ of being absorbed. These values for the probabilities lead to the correct classical limit for a beam containing a large number of photons.

Note that we have only been able to preserve the individuality of photons, in all cases, by abandoning the determinacy of classical theory, and adopting a fundamentally probabilistic approach. We have no way of knowing whether a given obliquely-polarized photon is going to be absorbed by, or transmitted through, the polarizing film. We only know the probability of each event occurring. This is a fairly sweeping statement. Recall, however, that the state of a photon is fully specified once its energy, direction of propagation, and polarization are known. If we imagine performing experiments using monochromatic light, normally incident on a polarizing film, with a particular oblique polarization, then the state of each individual photon in the beam is completely specified, and nothing remains to uniquely determine whether the photon is transmitted or absorbed by the film.

The previous discussion about the possible results of an experiment with a single obliquely polarized photon incident on a polarizing film answers all that can be legitimately asked about what happens to the photon when it reaches the film. Questions as to what determines whether the photon is transmitted or not, or how it changes its direction of polarization, are illegitimate, because they do not relate to the outcome of a possible experiment. Nevertheless, some further description is needed, in order to allow the results of this experiment to be correlated with the results of other experiments that can be performed using photons.

The further description provided by quantum mechanics is as follows. It is supposed that a photon polarized obliquely to the optic axis can be regarded as being partly in a state of polarization parallel to the axis, and partly in a state of polarization perpendicular to the axis. In other words, the oblique polarization state is some sort of superposition of two states of parallel and perpendicular polarization. Because there is nothing special about the orientation of the optic axis in our experiment, we deduce that any photon polarization state can be regarded as a superposition of two mutually perpendicular polarization states. (Recall, from Section 2.4.4, that there are only two independent polarizations of an electromagnetic wave.) When we cause a photon to encounter a polarizing film, we are subjecting it to an observation. In fact, we are observing whether it is polarized parallel or perpendicular to the film's optic axis. The effect of making this observation is to force the photon entirely into a state of parallel or perpendicular polarization. In other words, the photon has to jump suddenly from being partly in each of these two states to being entirely in one or the other of them. Which of the two states it will jump into cannot be predicted, but is governed by probability laws. If the photon jumps into a state of parallel polarization then it is absorbed. Otherwise, it is transmitted. Note that, in this example, the introduction of indeterminacy into the problem is clearly connected with the act of observation. In other words, the indeterminacy is related to the inevitable disturbance of the system associated with the act of observation.

4.1.5 Double-Slit Interference of Light

As was first described by Thomas Young in 1801, if a monochromatic light source illuminates a plate pierced by two parallel slits, and the light passing through the slits is observed on a screen located behind the plate, then bright and dark bands appear on the screen. This experiment is usually thought of as a demonstration that light is a wavelike phenomenon. In fact, the conventional explanation is that incident light waves pass through both slits, and then travel slightly different distances to a given point on the screen, where they interfere with one another to produce a bright band if the path difference is an integer multiple of a wavelength, and a dark band if the path difference is a half-integer multiple of a wavelength.

How do we account for double-slit interference at the individual photon level? In fact, in 1909, Geoffrey I. Taylor showed that an interference pattern is generated in a double-slit experiment even when the incident light intensity is so low that only a single photon could be in the apparatus at a given time. The only way in which to account for this result is to assume that an individual photon incident on the apparatus passes through both slits, and then interferes with itself when it reaches the screen. In other words, a photon in the apparatus is partly in a state in which it passed through one slit, and partly in a state in which it passed through the other. Moreover, the interference between these two states at the screen can only determine the probability of the photon being observed at a

given point on the screen.

4.1.6 Electron Diffraction

In 1927, George P. Thomson discovered that if a beam of electrons is made to pass through a thin gold foil then the regular atomic array in the foil acts as a sort of diffraction grating, so that, when a photographic film placed behind the foil is developed, an interference pattern is discernible. Independently, Clinton Davisson and Lester Germer found that electrons scattered by the surface of a nickel metal crystal display a diffraction pattern. Both these experimental results imply that electrons have wave-like properties. The electron wavelength, λ , or, alternatively, the wavenumber, $k = 2\pi/\lambda$, can be deduced from the spacing of the maxima in the interference pattern. Thomson, Davisson, and Germer found that the momentum, \mathbf{p} , of an electron is related to its wavevector, \mathbf{k} , according to the following simple relation:

$$\mathbf{p} = \hbar \mathbf{k}. \quad (4.5)$$

The associated wavelength, $\lambda = 2\pi/k$, is known as the *de Broglie wavelength*, because the previous relation was first hypothesized by Louis de Broglie in 1926. (See Section 4.1.9.)

It turns out that wave-particle duality only manifests itself on lengthscales less than, or of order, the de Broglie wavelength. Under normal circumstances, this wavelength is fairly small. For instance, the de Broglie wavelength of an electron is

$$\lambda_e = 1.2 \times 10^{-9} [E(\text{eV})]^{-1/2} \text{ m}, \quad (4.6)$$

where the electron energy is conveniently measured in units of electron-volts (eV). (An electron accelerated from rest through a potential difference of 1000 V acquires an energy of 1000 eV, and so on. Electrons in atoms typically have energies in the range 10 to 100 eV.)

4.1.7 Helium Diffraction

In 1930, Immanuel Estermann and Otto Stern obtained a diffraction pattern from a beam of room temperature helium atoms scattered off a lithium fluoride crystal. Estermann and Stern were able to demonstrate that Equation (4.5) applies to helium atoms—and, by implication, to protons and neutrons—as well as to electrons. The de Broglie wavelength of a proton is

$$\lambda_p = 2.9 \times 10^{-11} [E(\text{eV})]^{-1/2} \text{ m}. \quad (4.7)$$

4.1.8 Two-Source Particle Interference

In 1961, Claus Jönsson performed a double-slit experiment with electrons and obtained the expected interference pattern. In 1974, Pier G. Merli, Gian F. Missiroli, and Giulio Pozzi performed a more advanced double-slit experiment with electrons, in which only one electron was in the apparatus at a given time, and also obtained the expected interference pattern. These experiments

demonstrate that the behavior of electrons in double-slit experiments is analogous to those of photons described in Section 4.1.5. The only way in which to account for the results of the Merli-Missiroli-Pozzi experiment is to assume that an individual electron incident on a double-slit apparatus passes through both slits, and then interferes with itself when it reaches the detection screen. In other words, an electron in the apparatus is partly in a state in which it passed through one slit, and partly in a state in which it passed through the other. Moreover, the interference between these two states at the screen can only determine the probability of the electron being observed at a given point on the screen.

4.1.9 de Broglie's Hypothesis

In 1926, Louis de Broglie hypothesized that massive particles have wave-like properties, and that the angular frequency, ω , and wavenumber, \mathbf{k} , of a particle wave is related to the energy, E , and momentum, \mathbf{p} , of its constituent particles according to

$$E = \hbar \omega, \quad (4.8)$$

$$\mathbf{p} = \hbar \mathbf{k}. \quad (4.9)$$

Obviously, these relations are the same as the relations between the angular frequency and wavenumber of an electromagnetic wave and the energy and momentum of its constituent photons. (See Section 4.1.2.) As discussed in Sections 4.1.6 and 4.1.7, relation (4.9) can be verified experimentally. On the other hand, relation (4.8) is adopted on the basis of an analogy drawn between massive particles and photons.

4.2 Wave Mechanics

4.2.1 Wavefunctions

The basic premise of *wave mechanics* is that a massive particle of energy E and linear momentum p , moving in the x -direction (say), can be represented by a one-dimensional wavefunction of the form

$$\psi(x, t) = \psi_0 e^{i(kx - \omega t)}, \quad (4.10)$$

where the complex amplitude, ψ_0 , is arbitrary, while the angular frequency, ω , and the wavenumber, k , are related to the particle energy, E , and momentum, p , via the fundamental relations

$$E = \hbar \omega, \quad (4.11)$$

$$p = \hbar k. \quad (4.12)$$

(See Section 4.1.9.)

The one-dimensional wavefunction (4.10) is the solution of a one-dimensional wave equation that determines how the wavefunction evolves in time. As described in the next section, we can guess the form of this wave equation by drawing an analogy with classical physics.

4.2.2 Schrödinger's Equation

A classical particle of mass m , moving in a one-dimensional potential $U(x)$, satisfies the energy conservation equation

$$E = K + U, \quad (4.13)$$

where

$$K = \frac{p^2}{2m} \quad (4.14)$$

is the particle's kinetic energy. (See Sections 1.3.2 and 1.3.5.) Hence,

$$E \psi = (K + U) \psi \quad (4.15)$$

is a valid, but not obviously useful, wave equation.

However, it follows from Equations (4.10) and (4.11) that

$$\frac{\partial \psi}{\partial t} = -i \omega \psi_0 e^{i(kx - \omega t)} = -i \frac{E}{\hbar} \psi, \quad (4.16)$$

which can be rearranged to give

$$E \psi = i \hbar \frac{\partial \psi}{\partial t}. \quad (4.17)$$

Likewise, from Equations (4.10) and (4.12),

$$\frac{\partial \psi}{\partial x} = i k \psi_0 e^{i(kx - \omega t)} = i \frac{p}{\hbar} \psi, \quad (4.18)$$

which can be rearranged to give

$$p \psi = -i \hbar \frac{\partial \psi}{\partial x}. \quad (4.19)$$

It immediately follows that

$$p^2 \psi = -\hbar^2 \frac{\partial^2 \psi}{\partial x^2}. \quad (4.20)$$

Hence,

$$K \psi = \frac{p^2}{2m} \psi = -\frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial x^2}. \quad (4.21)$$

Thus, combining Equations (4.15), (4.17), and (4.21), we obtain

$$i \hbar \frac{\partial \psi}{\partial t} = -\frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial x^2} + U(x) \psi. \quad (4.22)$$

This equation, which is known as *Schrödinger's equation*—because it was first formulated by Erwin Schrödinger in 1926—is the fundamental equation of wave mechanics.

For a massive particle moving in free space (i.e., $U = 0$), the complex wavefunction (4.10) is a solution of Schrödinger's equation, (4.22), provided

$$\omega = \frac{\hbar}{2m} k^2. \quad (4.23)$$

The previous expression can be thought of as the dispersion relation for matter waves in free space. (See Section 4.2.6.) The associated *phase velocity* (i.e., propagation speed of a wave maximum) is

$$v_p = \frac{\omega}{k} = \frac{\hbar k}{2m} = \frac{p}{2m}, \quad (4.24)$$

where use has been made of Equation (4.12). However, this phase velocity is only half the classical velocity, $v = p/m$, of a massive (non-relativistic) particle.

4.2.3 Probability Interpretation of Wavefunction

After many false starts, physicists in the early 20th century eventually came to the conclusion that the only physical interpretation of a particle wavefunction that is consistent with experimental observations is probabilistic in nature. To be more exact, if $\psi(x, t)$ is the complex wavefunction of a given particle, moving in one dimension along the x -axis, then the probability of finding the particle between x and $x + dx$ at time t is

$$P(x, t) = |\psi(x, t)|^2 dx. \quad (4.25)$$

A probability is a real number lying in the range 0 to 1. An event that has a probability 0 is impossible. On the other hand, an event that has a probability 1 is certain to occur. An event that has a probability 1/2 (say) is such that in a very large number of identical trials the event occurs in half of the trials. (See Section 5.1.1.)

We can interpret

$$P(t) = \int_{-\infty}^{\infty} |\psi(x, t)|^2 dx \quad (4.26)$$

as the probability of the particle being found anywhere between $x = -\infty$ and $x = +\infty$ at time t . This follows, via induction, from the fundamental result in probability theory that the probability of the occurrence of one or other of two mutually exclusive events (such as the particle being found in two non-overlapping regions) is the sum (or integral) of the probabilities of the individual events. (For example, the probability of throwing a 1 on a six-sided die is 1/6. Likewise, the probability of throwing a 2 is 1/6. Hence, the probability of throwing a 1 or a 2 is $1/6 + 1/6 = 1/3$.) Assuming that the particle exists, it is certain that it will be found somewhere between $x = -\infty$ and $x = +\infty$ at time t . Because a certain event has probability 1, our probability interpretation of the wavefunction is only tenable provided

$$\int_{-\infty}^{\infty} |\psi(x, t)|^2 dx = 1 \quad (4.27)$$

at all times. A wavefunction that satisfies the previous condition—which is known as the *normalization condition*—is said to be properly normalized.

Suppose that we have a wavefunction, $\psi(x, t)$, which is such that it satisfies the normalization condition (4.27) at time $t = 0$. Furthermore, let the wavefunction evolve in time according to Schrödinger's equation, (4.22). Our probability interpretation of the wavefunction only makes sense if the normalization condition remains satisfied at all subsequent times. This follows because if the particle is certain to be found somewhere on the x -axis (which is the interpretation put on

the normalization condition) at time $t = 0$ then it is equally certain to be found somewhere on the x -axis at a later time (because we are not considering any physical process by which particles can be created or destroyed). Thus, it is necessary for us to demonstrate that Schrödinger's equation preserves the normalization of the wavefunction.

Taking Schrödinger's equation, and multiplying it by ψ^* (the complex conjugate of the wavefunction), we obtain

$$i\hbar \frac{\partial \psi}{\partial t} \psi^* = -\frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial x^2} \psi^* + U(x) |\psi|^2. \quad (4.28)$$

The complex conjugate of the previous expression yields

$$-i\hbar \frac{\partial \psi^*}{\partial t} \psi = -\frac{\hbar^2}{2m} \frac{\partial^2 \psi^*}{\partial x^2} \psi + U(x) |\psi|^2. \quad (4.29)$$

Here, use has been made of the readily demonstrated results $(\psi^*)^* = \psi$ and $i^* = -i$, as well as the fact that $U(x)$ is real. Taking the difference between the previous two expressions, we obtain

$$i\hbar \left(\frac{\partial \psi}{\partial t} \psi^* + \frac{\partial \psi^*}{\partial t} \psi \right) = -\frac{\hbar^2}{2m} \left(\frac{\partial^2 \psi}{\partial x^2} \psi^* - \frac{\partial^2 \psi^*}{\partial x^2} \psi \right), \quad (4.30)$$

which can be written

$$i\hbar \frac{\partial |\psi|^2}{\partial t} = -\frac{\hbar^2}{2m} \frac{\partial}{\partial x} \left(\frac{\partial \psi}{\partial x} \psi^* - \frac{\partial \psi^*}{\partial x} \psi \right). \quad (4.31)$$

Integrating in x , we get

$$i\hbar \frac{d}{dt} \int_{-\infty}^{\infty} |\psi|^2 dx = -\frac{\hbar^2}{2m} \left[\frac{\partial \psi}{\partial x} \psi^* - \frac{\partial \psi^*}{\partial x} \psi \right]_{-\infty}^{\infty}. \quad (4.32)$$

Finally, assuming that the wavefunction is localized in space; that is,

$$|\psi(x, t)| \rightarrow 0 \quad \text{as} \quad |x| \rightarrow \infty, \quad (4.33)$$

we obtain

$$\frac{d}{dt} \int_{-\infty}^{\infty} |\psi|^2 dx = 0. \quad (4.34)$$

It follows, from the preceding analysis, that if a localized wavefunction is properly normalized at $t = 0$ (i.e., if $\int_{-\infty}^{\infty} |\psi(x, 0)|^2 dx = 1$) then it will remain properly normalized as it evolves in time according to Schrödinger's equation.

A wavefunction that is not localized cannot be properly normalized, because its normalization integral $\int_{-\infty}^{\infty} |\psi|^2 dx$ is necessarily infinite. For such a wavefunction, $|\psi(x, t)|^2 dx$ gives the relative, rather than the absolute, probability of finding the particle between x and $x + dx$ at time t . In other words, [cf., Equation (4.25)]

$$P(x, t) \propto |\psi(x, t)|^2 dx. \quad (4.35)$$

4.2.4 Wave Packets

As we have seen, the wavefunction of a massive particle of momentum p and energy E , moving in free space along the x -axis, can be written

$$\psi(x, t) = \bar{\psi} e^{i(kx - \omega t)}, \quad (4.36)$$

where $k = p/\hbar$, $\omega = E/\hbar$, and $\bar{\psi}$ is a complex constant. Here, ω and k are linked via the matter-wave dispersion relation (4.23). Expression (4.36) represents a plane wave that propagates in the x -direction with the phase velocity $v_p = \omega/k$. However, it follows from Equation (4.24) that this phase velocity is only half of the classical velocity of a massive particle.

According to the discussion in the previous section, the most reasonable physical interpretation of the wavefunction is that $|\psi(x, t)|^2 dx$ is proportional to (assuming that the wavefunction is not properly normalized) the probability of finding the particle between x and $x+dx$ at time t . However, the modulus squared of the wavefunction (4.36) is $|\bar{\psi}|^2$, which is a constant that depends on neither x nor t . In other words, the previous wavefunction represents a particle that is equally likely to be found anywhere on the x -axis at all times. Hence, the fact that this wavefunction propagates at a phase velocity that does not correspond to the classical particle velocity has no observable consequences.

How can we write the wavefunction of a particle that is localized in x ? In other words, a particle that is more likely to be found at some positions on the x -axis than at others. It turns out that we can achieve this goal by forming a linear combination of plane waves of different wavenumbers; that is,

$$\psi(x, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \bar{\psi}(k) e^{i(kx - \omega t)} dk. \quad (4.37)$$

Here, $\bar{\psi}(k)$ represents the complex amplitude of plane waves of wavenumber k within this combination. In writing the previous expression, we are relying on the assumption that matter waves are superposable. In other words, it is possible to add two valid wave solutions to form a third valid wave solution. The ultimate justification for this assumption is that matter waves satisfy the linear wave equation (4.22).

There is a fundamental mathematical theorem, known as *Fourier's theorem*, that states that if

$$f(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \bar{f}(k) e^{ikx} dk, \quad (4.38)$$

then

$$\bar{f}(k) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) e^{-ikx} dx. \quad (4.39)$$

Here, $\bar{f}(k)$ is known as the Fourier transform of the function $f(x)$. We can use Fourier's theorem to find the k -space function $\bar{\psi}(k)$ that generates any given x -space wavefunction $\psi(x)$ at a given time.

For instance, suppose that at $t = 0$ the wavefunction of our particle takes the form

$$\psi(x, 0) = \frac{1}{[2\pi (\Delta x)^2]^{1/4}} \exp \left[ik_0 x - \frac{(x - x_0)^2}{4 (\Delta x)^2} \right]. \quad (4.40)$$

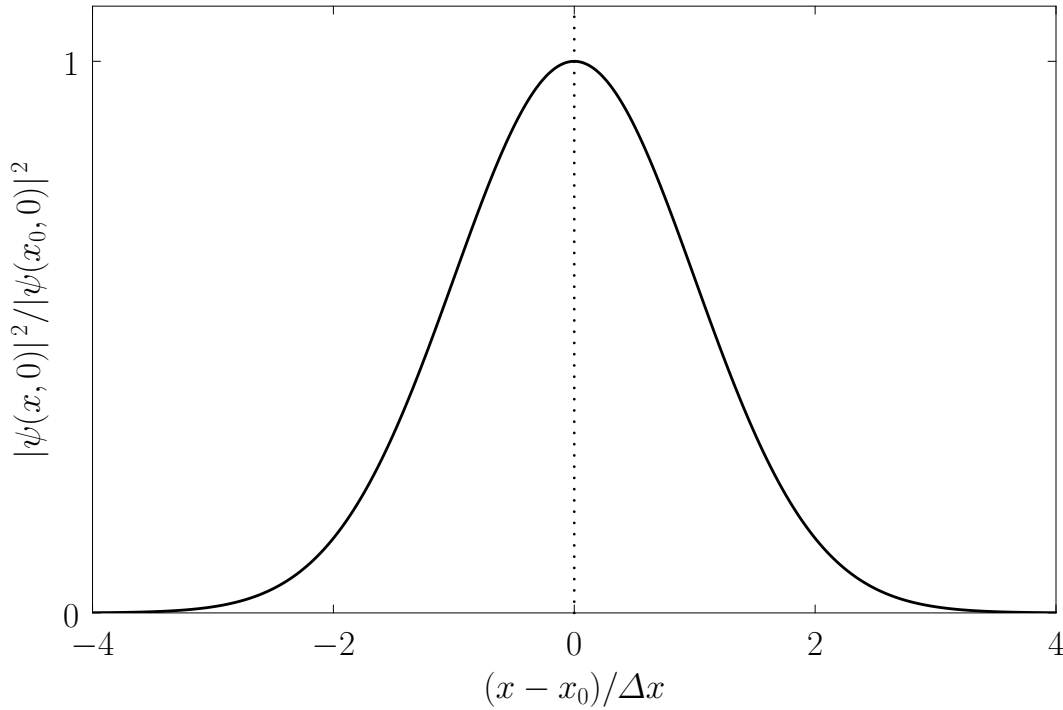


Figure 4.2: A one-dimensional Gaussian probability distribution.

Thus, the initial probability distribution for the particle's x -coordinate is

$$|\psi(x, 0)|^2 = \frac{1}{[2\pi(\Delta x)^2]^{1/2}} \exp\left[-\frac{(x - x_0)^2}{2(\Delta x)^2}\right]. \quad (4.41)$$

This particular distribution is called a *Gaussian distribution* (see Section 5.1.7), and is plotted in Figure 4.2. It can be seen that a measurement of the particle's position is most likely to yield the value x_0 , and very unlikely to yield a value which differs from x_0 by more than $3\Delta x$. Thus, Equation (4.40) is the wavefunction of a particle that is initially localized in some region of x -space, centered on $x = x_0$, whose width is of order Δx . This type of wavefunction is known as a *wave packet*.

It is easily demonstrated that the wavefunction (4.40) is properly normalized. In fact,

$$\begin{aligned} \int_{-\infty}^{\infty} |\psi(x, 0)|^2 dx &= \frac{1}{[2\pi(\Delta x)^2]^{1/2}} \int_{-\infty}^{\infty} \exp\left[-\frac{(x - x_0)^2}{2(\Delta x)^2}\right] dx \\ &= \frac{2^{1/2}(\Delta x)}{[2\pi(\Delta x)^2]^{1/2}} \int_{-\infty}^{\infty} e^{-y^2} dy = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} e^{-y^2} dy = 1. \end{aligned} \quad (4.42)$$

Here, $y = (x - x_0)/(2^{1/2}\Delta x)$, and use has been made of the standard result

$$\int_{-\infty}^{\infty} e^{-y^2} dy = \pi^{1/2}. \quad (4.43)$$

According to Equation (4.37),

$$\psi(x, 0) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \bar{\psi}(k) e^{ikx} dk. \quad (4.44)$$

Hence, we can employ Fourier's theorem to invert this expression to give

$$\bar{\psi}(k) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \psi(x, 0) e^{-ikx} dx. \quad (4.45)$$

Making use of Equation (4.40), we obtain

$$\begin{aligned} \bar{\psi}(k) &= \frac{1}{(2\pi)^{3/4} (\Delta x)^{1/2}} \int_{-\infty}^{\infty} \exp \left[-i(k - k_0)x - \frac{(x - x_0)^2}{4(\Delta x)^2} \right] dx \\ &= \frac{1}{(2\pi)^{3/4} (\Delta x)^{1/2}} \exp \left[-i(k - k_0)x_0 - \frac{(k - k_0)^2}{4(\Delta k)^2} \right] \int_{-\infty}^{\infty} \exp \left\{ - \left[\frac{x - x_0}{2\Delta x} + i\Delta x(k - k_0) \right]^2 \right\} dx \\ &= \frac{2\Delta x}{(2\pi)^{3/4} (\Delta x)^{1/2}} \exp \left[-i(k - k_0)x_0 - \frac{(k - k_0)^2}{4(\Delta k)^2} \right] \int_{-\infty}^{\infty} e^{-y^2} dy \\ &= \frac{1}{[2\pi(\Delta k)^2]^{1/4}} \exp \left[-i(k - k_0)x_0 - \frac{(k - k_0)^2}{4(\Delta k)^2} \right], \end{aligned} \quad (4.46)$$

where $y = (x - x_0)/(2\Delta x) + i\Delta x(k - k_0)$,

$$\Delta k = \frac{1}{2\Delta x}, \quad (4.47)$$

and use has been made of Equation (4.43).

If $|\psi(x, 0)|^2 dx$ is the probability that a measurement of the particle's position yields a value in the range x to $x + dx$ at time $t = 0$ then it stands to reason that $|\bar{\psi}(k)|^2 dk$ is the probability that a measurement of the particle's wavenumber yields a value in the range k to $k + dk$. (Recall that $p = \hbar k$, so a measurement of the particle's wavenumber, k , is equivalent to a measurement of the particle's momentum, p .) According to Equation (4.46),

$$|\bar{\psi}(k)|^2 = \frac{1}{[2\pi(\Delta k)^2]^{1/2}} \exp \left[-\frac{(k - k_0)^2}{2(\Delta k)^2} \right]. \quad (4.48)$$

This probability distribution is a Gaussian in k -space. [See Equation (4.41) and Figure 4.2.] Hence, a measurement of k is most likely to yield the value k_0 , and very unlikely to yield a value that differs from k_0 by more than $3\Delta k$. Note that the probability distribution (4.48) is properly normalized; that is, $\int_{-\infty}^{\infty} |\bar{\psi}(k)|^2 dk = 1$.

We have just seen that a wave packet with a Gaussian probability distribution of characteristic width Δx in x -space [see Equation (4.41)] is equivalent to a wave packet with a Gaussian probability distribution of characteristic width Δk in k -space [see Equation (4.48)], where

$$\Delta x \Delta k = \frac{1}{2}. \quad (4.49)$$

This illustrates an important property of wave packets. Namely, in order to construct a packet that is highly localized in x -space (i.e., with small Δx) we need to combine plane waves with a very wide range of different k -values (i.e., with large Δk). Conversely, if we only combine plane waves whose wavenumbers differ by a small amount (i.e., if Δk is small) then the resulting wave packet is highly extended in x -space (i.e., Δx is large).

4.2.5 Group Velocity

We have seen that Equation (4.40) is the wavefunction of a particle whose most probable position at time $t = 0$ is $x = x_0$. According to Equations (4.37) and (4.46), the wavefunction evolves in time as

$$\psi(x, t) = \frac{1}{(2\pi)^{3/4} (\Delta k)^{1/2}} \int_{-\infty}^{\infty} \exp \left[i k x - i \omega t - i (k - k_0) x_0 - \frac{(k - k_0)^2}{4 (\Delta k)^2} \right] dk. \quad (4.50)$$

Here, ω is related to k via the dispersion relation (4.23); in other words, $\omega = \omega(k)$. Now, the integrand on the right-hand side of the previous expression is strongly peaked at $k = k_0$. It follows that the only significant contribution to the corresponding integral comes from a small region of k -space centered on $k = k_0$. Let us Taylor expand the dispersion relation, $\omega = \omega(k)$, about $k = k_0$. Neglecting second-order terms in the expansion, we obtain

$$\omega \simeq \omega_0 + v_g (k - k_0), \quad (4.51)$$

where

$$\omega_0 = \omega(k_0), \quad (4.52)$$

$$v_g = \frac{d\omega(k_0)}{dk}. \quad (4.53)$$

Thus, we get

$$\begin{aligned} \psi(x, t) &= \frac{1}{(2\pi)^{3/4} (\Delta k)^{1/2}} \exp[i(k_0 x - \omega_0 t)] \int_{-\infty}^{\infty} \left[i (k - k_0) (x - x_0 - v_g t) - \frac{(k - k_0)^2}{4 (\Delta k)^2} \right] dk \\ &= \frac{2 \Delta k}{(2\pi)^{3/4} (\Delta k)^{1/2}} \exp \left[i (k_0 x - \omega_0 t) - \frac{(x - x_0 - v_g t)^2}{4 (\Delta x)^2} \right] \int_{-\infty}^{\infty} e^{-y^2} dy, \end{aligned} \quad (4.54)$$

where $y = (k - k_0)/(2 \Delta k) - i \Delta k (x - x_0 - v_g t)$, and use has been made of Equation (4.49). The previous equation reduces to

$$\psi(x, t) = \frac{1}{[2\pi (\Delta x)^2]^{1/4}} \exp \left[i (k_0 x - \omega_0 t) - \frac{(x - x_0 - v_g t)^2}{4 (\Delta x)^2} \right], \quad (4.55)$$

where use has been made of Equations (4.43) and (4.49). Hence, the probability of finding the particle between x and $x + dx$ at time t is $|\psi(x, t)|^2 dx$, where

$$|\psi(x, t)|^2 = \frac{1}{[2\pi (\Delta x)^2]^{1/2}} \exp \left[-\frac{(x - x_0 - v_g t)^2}{2 (\Delta x)^2} \right]. \quad (4.56)$$

It can be seen that the particle's most probable location at time t is

$$x = x_0 + v_g t. \quad (4.57)$$

If, as seems reasonable, we identify the velocity of the particle with the velocity of its most probable location then we deduce that the particle effectively moves at the so-called *group velocity*,

$$v_g = \frac{d\omega}{dk}, \quad (4.58)$$

rather than the phase velocity,

$$v_p = \frac{\omega}{k}. \quad (4.59)$$

Incidentally, the distinction between these two velocities is as follows. The phase velocity is the propagation velocity of an individual wave maximum, whereas the group velocity is the propagation velocity of an interference peak.

We have seen that a spatially localized particle moves at the group velocity, (4.58), rather than the phase velocity, (4.59). Making use of the matter-wave dispersion relation, (4.23), the group velocity is

$$v_g = \frac{\hbar k}{m} = \frac{p}{m}, \quad (4.60)$$

where use has been made of Equation (4.12). This velocity is identical to the classical velocity of a (non-relativistic) massive particle. We conclude that the matter-wave dispersion relation (4.23) is perfectly consistent with classical physics, as long as we recognize that particles must be identified with wave packets (which propagate at the group velocity) rather than plane waves (which propagate at the phase velocity).

4.2.6 Wave Dispersion

Equation (4.56) indicates that, as a wave packet propagates, its envelope remains the same shape. Actually, this result is misleading, and is only obtained because of the neglect of second-order terms in the expansion (4.51). If we keep more terms in this expansion then we can show that the wave packet does actually change shape as it propagates. However, this demonstration is most readily effected by means of the following simple argument. The packet extends in Fourier space from $k_0 - \Delta k/2$ to $k_0 + \Delta k/2$. Thus, part of the packet propagates at the velocity $v_g(k_0 - \Delta k/2)$, and part at the velocity $v_g(k_0 + \Delta k/2)$. Consequently, the packet spreads out as it propagates, because some parts of it move faster than others. Roughly speaking, the spatial extent of the packet in real space grows as

$$\Delta x \sim (\Delta x)_0 + [v_g(k_0 + \Delta k/2) - v_g(k_0 - \Delta k/2)] t \sim (\Delta x)_0 + \frac{dv_g(k_0)}{dk} \Delta k t, \quad (4.61)$$

where $(\Delta x)_0 \sim (\Delta k)^{-1}$ is the extent of the packet at $t = 0$. Hence, from Equation (4.53),

$$\Delta x \sim (\Delta x)_0 + \frac{d^2\omega(k_0)}{dk^2} \frac{t}{(\Delta x)_0}. \quad (4.62)$$

We, thus, conclude that the spatial extent of the packet grows linearly in time, at a rate proportional to the second derivative of $\omega(k)$ with respect to k (evaluated at the packet's central wavenumber). This effect is known as *wave dispersion*. Furthermore, it is clear that the relation $\omega = \omega(k)$ governs the degree of wave dispersion, which explains why it is called a dispersion relation.

Note that electromagnetic wave packets, which are governed by the linear dispersion relation $\omega = kc$, do not disperse as they propagate, because $d^2\omega/dk^2 = 0$. Particle wave packets, on the other hand, are governed by the quadratic dispersion relation (4.23) and, therefore, disperse as they propagate. In fact, it follows from Equations (4.23) and (4.62) that the width of a particle wave packet grows in time as

$$\Delta x \simeq (\Delta x)_0 + \frac{\hbar}{m} \frac{t}{(\Delta x)_0}. \quad (4.63)$$

For example, if an electron wave packet is initially localized in a region of atomic dimensions (i.e., $\Delta x \sim 10^{-10}$ m) then the width of the packet doubles in about 10^{-16} s.

4.2.7 Heisenberg's Uncertainty Principle

According to the analysis contained in Section 4.2.4, a particle wave packet that is initially localized in x -space, with characteristic width Δx , is also localized in k -space, with characteristic width $\Delta k = 1/(2\Delta x)$. However, as time progresses, the width of the wave packet in x -space increases [see Equation (4.63)], while that of the packet in k -space stays the same [because $\bar{\psi}(k)$ is given by Equation (4.45) at all times]. Hence, in general, we can say that

$$\Delta x \Delta k \gtrsim \frac{1}{2}. \quad (4.64)$$

Furthermore, we can interpret Δx and Δk as characterizing our uncertainty regarding the values of the particle's position and wavenumber, respectively.

A measurement of a particle's wavenumber, k , is equivalent to a measurement of its momentum, p , because $p = \hbar k$. Hence, an uncertainty in k of order Δk translates to an uncertainty in p of order $\Delta p = \hbar \Delta k$. It follows, from the previous inequality, that

$$\Delta x \Delta p \gtrsim \frac{\hbar}{2}. \quad (4.65)$$

This result is known as the *Heisenberg uncertainty principle*, and was first proposed by Werner Heisenberg in 1927. According to this principle, it is impossible to simultaneously measure the position and momentum of a particle (exactly). Indeed, a good knowledge of the particle's position implies a poor knowledge of its momentum, and vice versa. The uncertainty principle is a direct consequence of representing particles as waves.

It is apparent, from Equation (4.63), that a particle wave packet of initial spatial extent $(\Delta x)_0$ spreads out in such a manner that its spatial extent becomes

$$\Delta x \sim \frac{\hbar t}{m(\Delta x)_0} \quad (4.66)$$

at large t . It is readily demonstrated that this spreading of the wave packet is a consequence of the uncertainty principle. Indeed, because the initial uncertainty in the particle's position is $(\Delta x)_0$, it follows that the uncertainty in its momentum is of order $\hbar/(\Delta x)_0$. This translates to an uncertainty in velocity of $\Delta v = \hbar/[m(\Delta x)_0]$. Thus, if we imagine that part of the wave packet propagates at $v_0 + \Delta v/2$, and another part at $v_0 - \Delta v/2$, where v_0 is the mean propagation velocity, then it follows that the wave packet will spread out as time progresses. Indeed, at large t , we expect the width of the wave packet to be

$$\Delta x \sim \Delta v t \sim \frac{\hbar t}{m(\Delta x)_0}, \quad (4.67)$$

which is identical to Equation (4.66). Evidently, the spreading of a particle wave packet, as time progresses, should be interpreted as representing an increase in our uncertainty regarding the particle's position, rather than an increase in the spatial extent of the particle itself.

4.2.8 Wavefunction Collapse

Consider a spatially extended wavefunction, $\psi(x, t)$. According to our standard interpretation, $|\psi(x, t)|^2 dx$ is proportional to the probability of a measurement of the particle's position yielding a value in the range x to $x + dx$ at time t . Thus, if the wavefunction is extended then there is a wide range of likely values that such a measurement could give. Suppose, however, that we make a measurement of the particle's position, and obtain the value x_0 . We now know that the particle is located at $x = x_0$. If we make another measurement, immediately after the first one, then what value would we expect to obtain? Common sense tells us that we should obtain the same value, x_0 , because the particle cannot have shifted position appreciably in an infinitesimal time interval. Thus, immediately after the first measurement, a measurement of the particle's position is certain to give the value x_0 , and has no chance of giving any other value. This implies that the wavefunction must have collapsed to some sort of "spike" function, centered on $x = x_0$. This idea is illustrated in Figure 4.3. As soon as the wavefunction collapses, it starts to expand again, as described in Section 4.2.6. Thus, the second measurement must be made reasonably quickly after the first one, otherwise the same result will not necessarily be obtained.

The preceding discussion illustrates an important point in wave mechanics. That is, the wavefunction of a massive particle changes discontinuously (in time) whenever a measurement of the particle's position is made. We conclude that there are two types of time evolution of the wavefunction in wave mechanics. First, there is a smooth evolution that is governed by Schrödinger's equation. This evolution takes place between measurements. Second, there is a discontinuous evolution that takes place each time a measurement is made.

4.2.9 Stationary States

Consider separable solutions to Schrödinger's equation of the form

$$\psi(x, t) = \psi(x) e^{-i\omega t}. \quad (4.68)$$

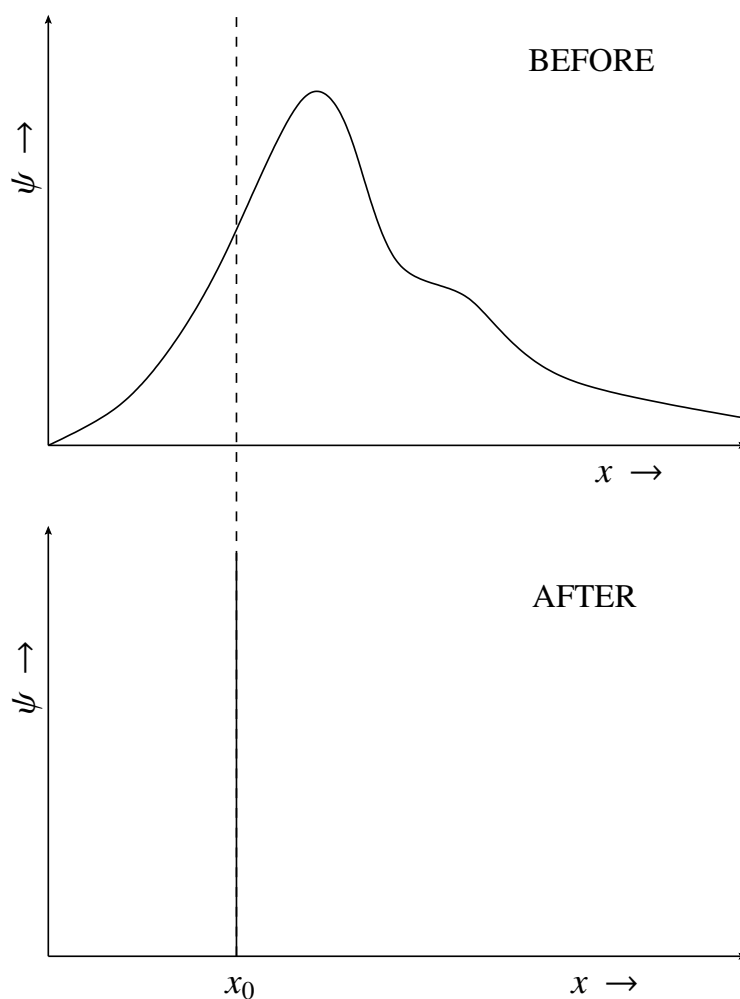


Figure 4.3: Collapse of the wavefunction upon measurement of x .

According to Equation (4.17), such solutions have definite energies $E = \hbar \omega$. For this reason, they are usually written

$$\psi(x, t) = \psi(x) e^{-i(E/\hbar)t}. \quad (4.69)$$

The probability of finding the particle between x and $x + dx$ at time t is

$$P(x, t) = |\psi(x, t)|^2 dx = |\psi(x)|^2 dx. \quad (4.70)$$

This probability is time independent. For this reason, states whose wavefunctions are of the form (4.69) are known as *stationary states*. Moreover, $\psi(x)$ is called a stationary wavefunction. Substituting (4.69) into Schrödinger's equation, (4.22), we obtain the following differential equation for $\psi(x)$;

$$-\frac{\hbar^2}{2m} \frac{d^2\psi}{dx^2} + U(x)\psi = E\psi. \quad (4.71)$$

This equation is called the *time-independent Schrödinger equation*.

4.3 One-Dimensional Wave Mechanics

4.3.1 Particle in Infinite Square Potential Well

Consider a particle trapped in a one-dimensional square potential well, of infinite depth, which is such that

$$U(x) = \begin{cases} 0 & 0 \leq x \leq a \\ \infty & \text{otherwise} \end{cases} . \quad (4.72)$$

The particle is excluded from the region $x < 0$ or $x > a$, so $\psi = 0$ in this region (i.e., there is zero probability of finding the particle outside the well). Within the well, a particle of definite energy E has a stationary wavefunction, $\psi(x)$, that satisfies

$$-\frac{\hbar^2}{2m} \frac{d^2\psi}{dx^2} = E\psi. \quad (4.73)$$

[See Equation (4.71).] The boundary conditions are

$$\psi(0) = \psi(a) = 0. \quad (4.74)$$

This follows because $\psi = 0$ in the region $x < 0$ or $x > a$, and $\psi(x)$ must be continuous [because a discontinuous wavefunction would generate a singular term (i.e., the term involving $d^2\psi/dx^2$) in the time-independent Schrödinger equation, (4.71), that could not be balanced, even by an infinite potential].

Let us search for solutions to Equation (4.73) of the form

$$\psi(x) = \psi_0 \sin(kx), \quad (4.75)$$

where ψ_0 is a constant. It follows that

$$\frac{\hbar^2 k^2}{2m} = E. \quad (4.76)$$

The solution (4.75) automatically satisfies the boundary condition $\psi(0) = 0$. The second boundary condition, $\psi(a) = 0$, leads to a quantization of the wavenumber; that is,

$$k = n \frac{\pi}{a}, \quad (4.77)$$

where $n = 1, 2, 3$, et cetera. (A “quantized” quantity is one that can only take certain discrete values.) According to Equation (4.76), the energy is also quantized. In fact, $E = E_n$, where

$$E_n = n^2 \frac{\hbar^2 \pi^2}{2m a^2}. \quad (4.78)$$

Thus, the allowed wavefunctions for a particle trapped in a one-dimensional square potential well of infinite depth are

$$\psi_n(x, t) = A_n \sin\left(n\pi \frac{x}{a}\right) \exp\left(-i n^2 \frac{E_1}{\hbar} t\right), \quad (4.79)$$

where n is a positive integer, and A_n a constant. We cannot have $n = 0$, because, in this case, we obtain a null wavefunction—that is, $\psi = 0$, everywhere—which corresponds to a non-existent state. Furthermore, if n takes a negative integer value then it generates exactly the same wavefunction as the corresponding positive integer value (assuming $A_{-n} = -A_n$).

The constant A_n , appearing in the previous wavefunction, can be determined from the constraint that the wavefunction be properly normalized. For the case under consideration, the normalization condition (4.27) reduces to

$$\int_0^a |\psi(x)|^2 dx = 1. \quad (4.80)$$

It follows from Equation (4.79) that $|A_n|^2 = 2/a$. Hence, the properly normalized version of the wavefunction (4.79) is

$$\psi_n(x, t) = \left(\frac{2}{a}\right)^{1/2} \sin\left(n\pi \frac{x}{a}\right) \exp\left(-i n^2 \frac{E_1}{\hbar} t\right). \quad (4.81)$$

Figure 4.4 shows the first four properly normalized stationary wavefunctions for a particle trapped in a one-dimensional square potential well of infinite depth; that is, $\psi_n(x) = (2/a)^{1/2} \sin(n\pi x/a)$, for $n = 1$ to 4.

At first sight, it seems rather strange that the lowest possible energy for a particle trapped in a one-dimensional potential well is not zero, as would be the case in classical mechanics, but rather $E_1 = \hbar^2 \pi^2 / (2ma^2)$. In fact, as explained in the following, this residual energy is a direct consequence of Heisenberg's uncertainty principle. A particle trapped in a one-dimensional well of width a is likely to be found anywhere inside the well. Thus, the uncertainty in the particle's position is $\Delta x \sim a$. It follows from the uncertainty principle, (4.65), that

$$\Delta p \gtrsim \frac{\hbar}{2\Delta x} \sim \frac{\hbar}{a}. \quad (4.82)$$

In other words, the particle cannot have zero momentum. In fact, the particle's momentum must be at least $p \sim \hbar/a$. However, for a free particle, $E = p^2/2m$. Hence, the residual energy associated with the particle's residual momentum is

$$E \sim \frac{p^2}{m} \sim \frac{\hbar^2}{ma^2} \sim E_1. \quad (4.83)$$

This type of residual energy, which often occurs in quantum mechanical systems, and has no equivalent in classical mechanics, is called *zero-point energy*.

The most general wavefunction for a particle trapped in a one-dimensional square potential well, of infinite depth, is a superposition of all of the possible stationary states. That is,

$$\psi(x, t) = \sum_{n=1, \infty} a_n \psi_n(x, t), \quad (4.84)$$

where the a_n are complex numbers, and the $\psi_n(x, t)$ are specified in Equation (4.81). Consider

$$\int_0^a |\psi(x, t)|^2 dx = \sum_{n, m=1, \infty} a_n a_m^* \int_0^a \psi_n(x, t) \psi_m^*(x, t) dx$$

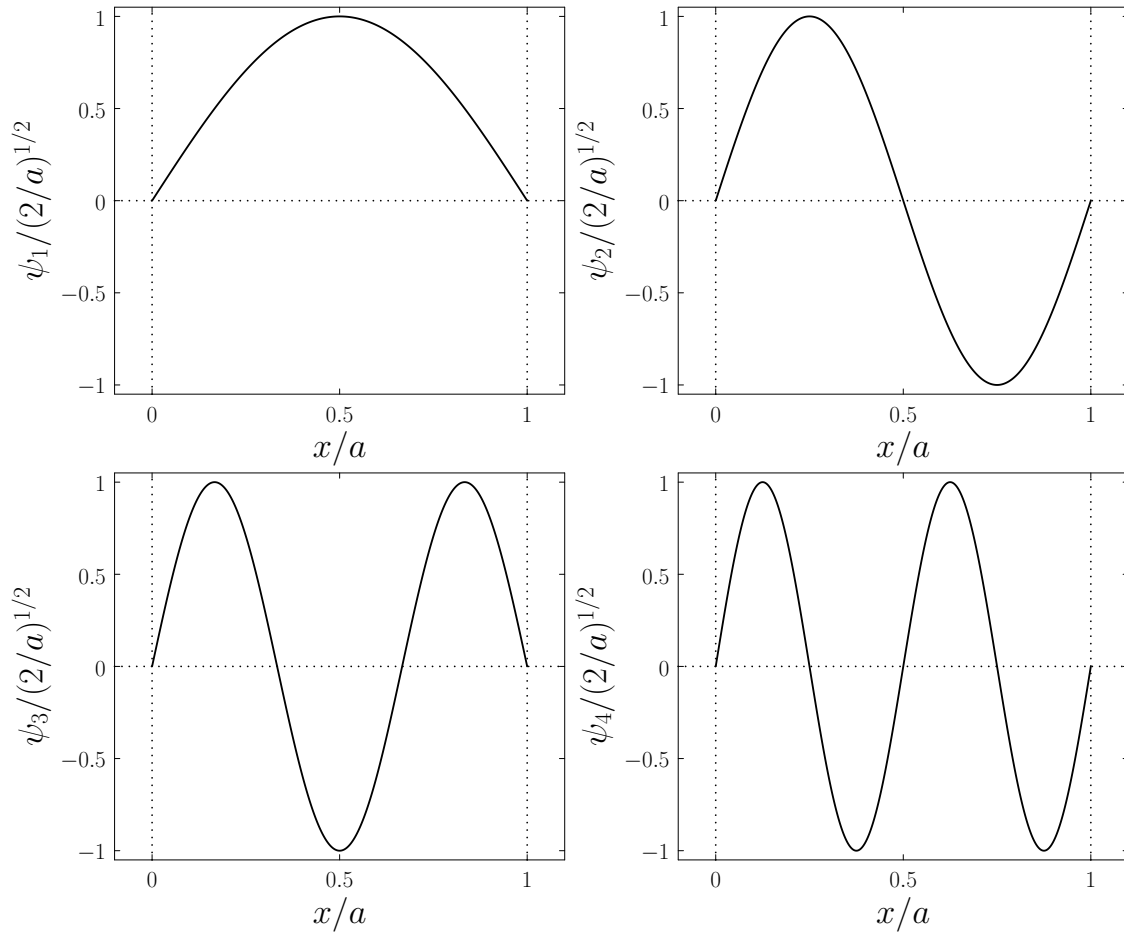


Figure 4.4: First four stationary wavefunctions for a particle trapped in a one-dimensional square potential well of infinite depth.

$$\begin{aligned}
 &= \sum_{n,m=1,\infty} a_n a_m^* \frac{2}{a} \int_0^a \sin\left(n\pi \frac{x}{a}\right) \sin\left(m\pi \frac{x}{a}\right) dx \\
 &= \sum_{n,m=1,\infty} a_n a_m^* \frac{2}{\pi} \int_0^\pi \sin(n\theta) \sin(m\theta) d\theta.
 \end{aligned} \tag{4.85}$$

However,

$$\frac{2}{\pi} \int_0^\pi \sin(n\theta) \sin(m\theta) d\theta = \delta_{nm}, \tag{4.86}$$

where δ_{nm} , which is known as a *Kronecker delta*, takes the value 1 if $n = m$, and 0 otherwise. Hence, we deduce that

$$\int_0^a |\psi(x, t)|^2 dx = \sum_{n=1,\infty} |a_n|^2. \tag{4.87}$$

Thus, the wavefunction (4.84) is properly normalized provided

$$\sum_{n=1,\infty} |a_n|^2 = 1. \quad (4.88)$$

Suppose that we make a measurement of the energy of a particle whose wavefunction is specified by Equation (4.84). Given that the wavefunction is a superposition of stationary states associated with the quantized energies E_n [see Equation (4.78)], it seems reasonable to assume that the measurement will result in one of these energies. In fact, according to quantum mechanics, the probability that a measurement of the particle's energy will give the result E_n is $|a_n|^2$. Thus, we can see that the normalization condition (4.88) ensures that the sum of all of these probabilities is unity. This must be the case, because a measurement of the particle's energy is certain to give one of the allowed energies. Suppose that we make a measurement of the particle's energy, and obtain the result E_n . A second measurement, made immediately after the first, must yield the same result. In other words, immediately after the first measurement, the particle's wavefunction must be such that a measurement of its energy is certain to give the result E_n , and has no chance of giving the result E_m , where $m \neq n$. This implies that $|a_n|^2 = 1$ and $|a_m|^2 = 0$, where $m \neq n$. We conclude that, after the first measurement, the particle's wavefunction is $\psi_n(x, t)$. This is another example of the collapse of a wavefunction consequent on a measurement. (See Section 4.2.8.)

4.3.2 Particle in Finite Square Potential Well

Consider a particle of mass m trapped in a one-dimensional, square, potential well of width a and finite depth $V > 0$. Suppose that the potential takes the form

$$U(x) = \begin{cases} -V & |x| \leq a/2 \\ 0 & \text{otherwise} \end{cases}. \quad (4.89)$$

Here, we have adopted the standard convention that $U(x) \rightarrow 0$ as $|x| \rightarrow \infty$. This convention is useful because, just as in classical mechanics, a particle whose overall energy, E , is negative is bound in the well (i.e., it cannot escape to infinity), whereas a particle whose overall energy is positive is unbound. (See Section 1.3.6.) Because we are interested in bound particles, we shall assume that $E < 0$. We shall also assume that $E + V > 0$, in order to allow the particle to have a positive kinetic energy inside the well.

Let us search for a stationary state

$$\psi(x, t) = \psi(x) e^{-i(E/\hbar)t}, \quad (4.90)$$

whose stationary wavefunction, $\psi(x)$, satisfies the time-independent Schrödinger equation, (4.71). Solutions to Equation (4.71) in the symmetric [i.e., $U(-x) = U(x)$] potential (4.89) are either totally symmetric [i.e., $\psi(-x) = \psi(x)$] or totally antisymmetric [i.e., $\psi(-x) = -\psi(x)$]. Moreover, the solutions must satisfy the boundary condition

$$\psi \rightarrow 0 \quad \text{as} \quad |x| \rightarrow \infty, \quad (4.91)$$

otherwise they would not correspond to bound states.

Let us, first of all, search for a totally-symmetric solution. In the region to the left of the well (i.e., $x < -a/2$), the solution of the time-independent Schrödinger equation that satisfies the boundary condition $\psi \rightarrow 0$ as $x \rightarrow -\infty$ is

$$\psi(x) = A e^{qx}, \quad (4.92)$$

where

$$q = \sqrt{\frac{2m(-E)}{\hbar^2}}, \quad (4.93)$$

and A is a constant. By symmetry, the solution in the region to the right of the well (i.e., $x > a/2$) is

$$\psi(x) = A e^{-qx}. \quad (4.94)$$

The solution inside the well (i.e., $|x| \leq a/2$) that satisfies the symmetry constraint $\psi(-x) = \psi(x)$ is

$$\psi(x) = B \cos(kx), \quad (4.95)$$

where

$$k = \sqrt{\frac{2m(V+E)}{\hbar^2}}, \quad (4.96)$$

and B is a constant. The appropriate matching conditions at the edges of the well (i.e., $x = \pm a/2$) are that $\psi(x)$ and $d\psi(x)/dx$ both be continuous [because a discontinuity in the wavefunction, or its first derivative, would generate a singular term in the time-independent Schrödinger equation (i.e., the term involving $d^2\psi/dx^2$) that could not be balanced]. The matching conditions yield

$$q = k \tan(ka/2). \quad (4.97)$$

Let $y = ka/2$. It follows that

$$E = E_0 y^2 - V, \quad (4.98)$$

where

$$E_0 = \frac{2\hbar^2}{ma^2}. \quad (4.99)$$

Moreover, Equation (4.97) becomes

$$\frac{\sqrt{\lambda - y^2}}{y} = \tan y, \quad (4.100)$$

with

$$\lambda = \frac{V}{E_0}. \quad (4.101)$$

Here, y must lie in the range $0 < y < \lambda^{1/2}$, in order to ensure that E lies in the range $-V < E < 0$.

The solutions of Equation (4.100) correspond to the intersection of the curve $(\lambda - y^2)^{1/2}/y$ with the curve $\tan y$. Figure 4.5 shows these two curves plotted for a particular value of λ . In this case,

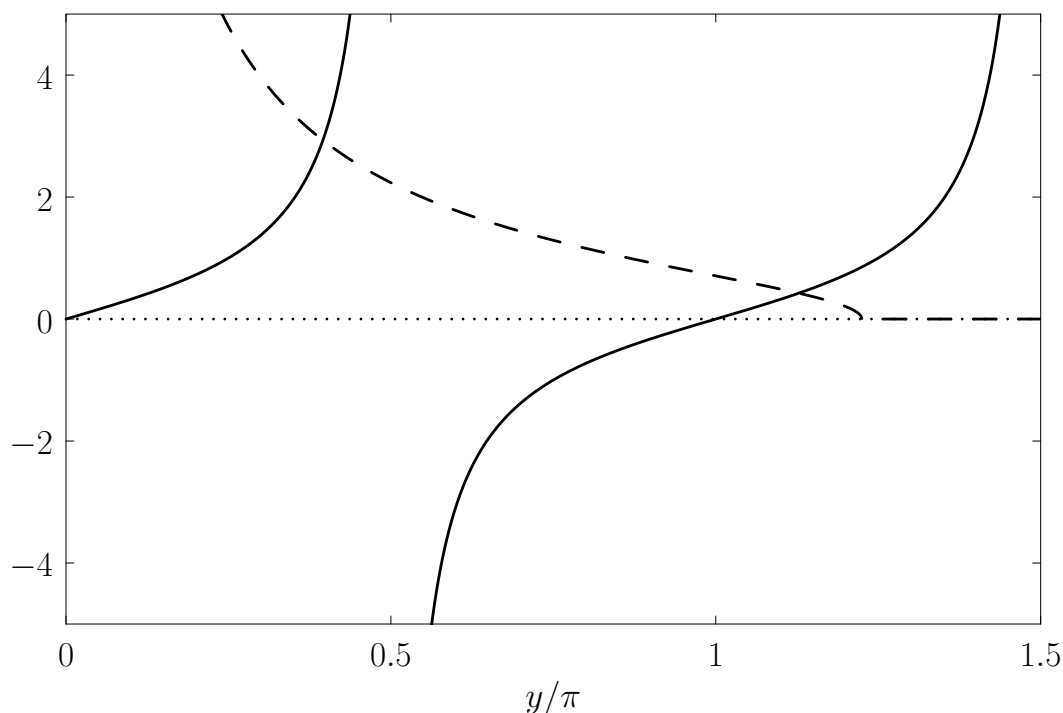


Figure 4.5: The curves $\tan y$ (solid) and $(\lambda - y^2)^{1/2}/y$ (dashed), calculated for $\lambda = 1.5\pi^2$. The latter curve takes the value 0 when $y > \lambda^{1/2}$.

the curves intersect twice, indicating the existence of two totally-symmetric bound states in the well. It is apparent, from the figure, that as λ increases (i.e., as the well becomes deeper) there are more and more bound states. However, it is also apparent that there is always at least one totally-symmetric bound state, no matter how small λ becomes (i.e., no matter how shallow the well becomes). In the limit $\lambda \gg 1$ (i.e., the limit in which the well is very deep), the solutions to Equation (4.100) asymptote to the roots of $\tan y = \infty$. This gives $y = (2n - 1)\pi/2$, where n is a positive integer, or

$$k = (2n - 1) \frac{\pi}{a}. \quad (4.102)$$

These solutions are equivalent to the odd- n infinite-depth potential well solutions specified by Equation (4.77).

For the case of a totally-antisymmetric bound state, similar analysis to the preceding yields

$$-\frac{y}{\sqrt{\lambda - y^2}} = \tan y. \quad (4.103)$$

The solutions of this equation correspond to the intersection of the curve $\tan y$ with the curve $-y/(\lambda - y^2)^{1/2}$. Figure 4.6 shows these two curves plotted for the same value of λ as that used in Figure 4.5. In this case, the curves intersect once, indicating the existence of a single totally-antisymmetric bound state in the well. It is, again, apparent, from the figure, that as λ increases (i.e., as the well becomes deeper) there are more and more bound states. However, it is also apparent

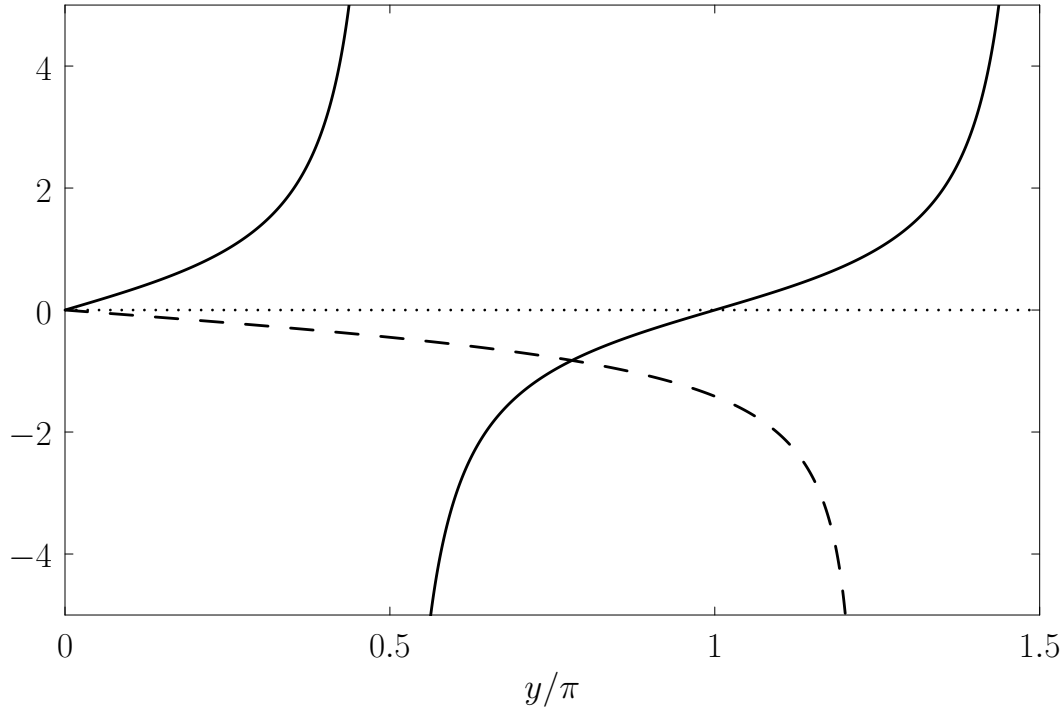


Figure 4.6: The curves $\tan y$ (solid) and $-y/(\lambda - y^2)^{1/2}$ (dashed), calculated for $\lambda = 1.5\pi^2$.

that when λ becomes sufficiently small [i.e., $\lambda < (\pi/2)^2$] then there is no totally antisymmetric bound state. In other words, a very shallow potential well always possesses a totally-symmetric bound state, but does not generally possess a totally-antisymmetric bound state. In the limit $\lambda \gg 1$ (i.e., the limit in which the well becomes very deep), the solutions to Equation (4.103) asymptote to the roots of $\tan y = 0$. This gives $y = n\pi$, where n is a positive integer, or

$$k = 2n \frac{\pi}{a}. \quad (4.104)$$

These solutions are equivalent to the even- n infinite-depth potential well solutions specified by Equation (4.77).

Probably the most surprising aspect of the bound states that we have just described is the possibility of finding the particle outside the well; that is, in the region $|x| > a/2$ where $U(x) > E$. This follows from Equation (4.94) and (4.95) because the ratio $A/B = \exp(qa/2) \cos(ka/2)$ is not necessarily zero. Such behavior is strictly forbidden in classical mechanics, according to which a particle of energy E is restricted to regions of space where $E > U(x)$. (See Section 1.3.6.) In fact, in the case of the ground state (i.e., the lowest energy symmetric state) it is possible to demonstrate that the probability of a measurement finding the particle outside the well is

$$P_{\text{out}} \simeq 1 - 2\lambda \quad (4.105)$$

for a shallow well (i.e., $\lambda \ll 1$), and

$$P_{\text{out}} \simeq \frac{\pi^2}{4} \frac{1}{\lambda^{3/2}} \quad (4.106)$$

for a deep well (i.e., $\lambda \gg 1$). It follows that the particle is very likely to be found outside a shallow well, and there is a small, but finite, probability of it being found outside a deep well. In fact, the probability of finding the particle outside the well only goes to zero in the case of an infinitely deep well (i.e., $\lambda \rightarrow \infty$).

4.3.3 Square Potential Barrier

Consider a particle of mass m and energy $E > 0$ interacting with the simple, one-dimensional, potential barrier

$$U(x) = \begin{cases} V & \text{for } 0 \leq x \leq a \\ 0 & \text{otherwise} \end{cases}, \quad (4.107)$$

where $V > 0$. In the regions to the left and to the right of the barrier, the stationary wavefunction, $\psi(x)$, satisfies

$$\frac{d^2\psi}{dx^2} = -k^2 \psi, \quad (4.108)$$

where

$$k = \sqrt{\frac{2mE}{\hbar^2}}. \quad (4.109)$$

Let us adopt the following solution of the previous equation to the left of the barrier (i.e., $x < 0$):

$$\psi(x) = e^{ikx} + R e^{-ikx}. \quad (4.110)$$

This solution consists of a plane wave of unit amplitude traveling to the right [because the full wavefunction is multiplied by a factor $\exp(-iEt/\hbar)$], and a plane wave of complex amplitude R traveling to the left. We interpret the first plane wave as an incident particle, and the second as a particle reflected by the potential barrier. Hence, $|R|^2$ is the probability of reflection.

Let us adopt the following solution to Equation (4.108) to the right of the barrier (i.e. $x > a$):

$$\psi(x) = T e^{ikx}. \quad (4.111)$$

This solution consists of a plane wave of complex amplitude T traveling to the right. We interpret the plane wave as a particle transmitted through the barrier. Hence, $|T|^2$ is the probability of transmission.

Let us, first of all, consider the situation in which $E > V$. In this case, according to classical mechanics, the particle slows down as it passes through the barrier, but is otherwise unaffected. In other words, the classical probability of reflection is zero, and the classical probability of transmission is unity. However, this is not necessarily the case in wave mechanics. In fact, inside the barrier (i.e., $0 \leq x \leq a$), $\psi(x)$ satisfies

$$\frac{d^2\psi}{dx^2} = -q^2 \psi, \quad (4.112)$$

where

$$q = \sqrt{\frac{2m(E - V)}{\hbar^2}}. \quad (4.113)$$

The general solution to Equation (4.112) takes the form

$$\psi(x) = A e^{iqx} + B e^{-iqx}. \quad (4.114)$$

Continuity of ψ and $d\psi/dx$ at the left edge of the barrier (i.e., $x = 0$) yields

$$1 + R = A + B, \quad (4.115)$$

$$k(1 - R) = q(A - B). \quad (4.116)$$

Likewise, continuity of ψ and $d\psi/dx$ at the right edge of the barrier (i.e., $x = a$) gives

$$A e^{iqa} + B e^{-iqa} = T e^{ika}, \quad (4.117)$$

$$q(A e^{iqa} - B e^{-iqa}) = kT e^{ika}. \quad (4.118)$$

After considerable algebra, the previous four equations yield

$$|T|^2 = 1 - |R|^2 = \frac{4k^2 q^2}{4k^2 q^2 + (k^2 - q^2)^2 \sin^2(qa)}. \quad (4.119)$$

The fact that $|R|^2 + |T|^2 = 1$ ensures that the probabilities of reflection and transmission sum to unity, as must be the case, because reflection and transmission are the only possible outcomes for a particle incident on the barrier.

The reflection and transmission probabilities obtained from Equation (4.119) are plotted in Figures 4.7 and 4.8. It can be seen, from Figure 4.7, that the classical result, $|R|^2 = 0$ and $|T|^2 = 1$, is obtained in the limit where the height of the barrier is relatively small (i.e., $V \ll E$). However, if V is of order E then there is a substantial probability that the incident particle will be reflected by the barrier. According to classical physics, reflection is impossible when $V < E$.

It can also be seen, from Figure 4.8, that at certain barrier widths the probability of reflection goes to zero. It turns out that this is true irrespective of the energy of the incident particle. It is evident, from Equation (4.119), that these special barrier widths correspond to

$$qa = n\pi, \quad (4.120)$$

where $n = 1, 2, 3, \dots$. In other words, the special barrier widths are integer multiples of half the de Broglie wavelength of the particle inside the barrier. There is no reflection at the special barrier widths because, at these widths, the backward traveling wave reflected from the left edge of the barrier interferes destructively with the similar wave reflected from the right edge of the barrier to give zero net reflected wave.

Let us now consider the situation in which $E < V$. In this case, according to classical mechanics, the particle is unable to penetrate the barrier, so the coefficient of reflection is unity, and the coefficient of transmission zero. However, this is not necessarily the case in wave mechanics. In fact, inside the barrier (i.e., $0 \leq x \leq a$), $\psi(x)$ satisfies

$$\frac{d^2\psi}{dx^2} = q^2 \psi, \quad (4.121)$$

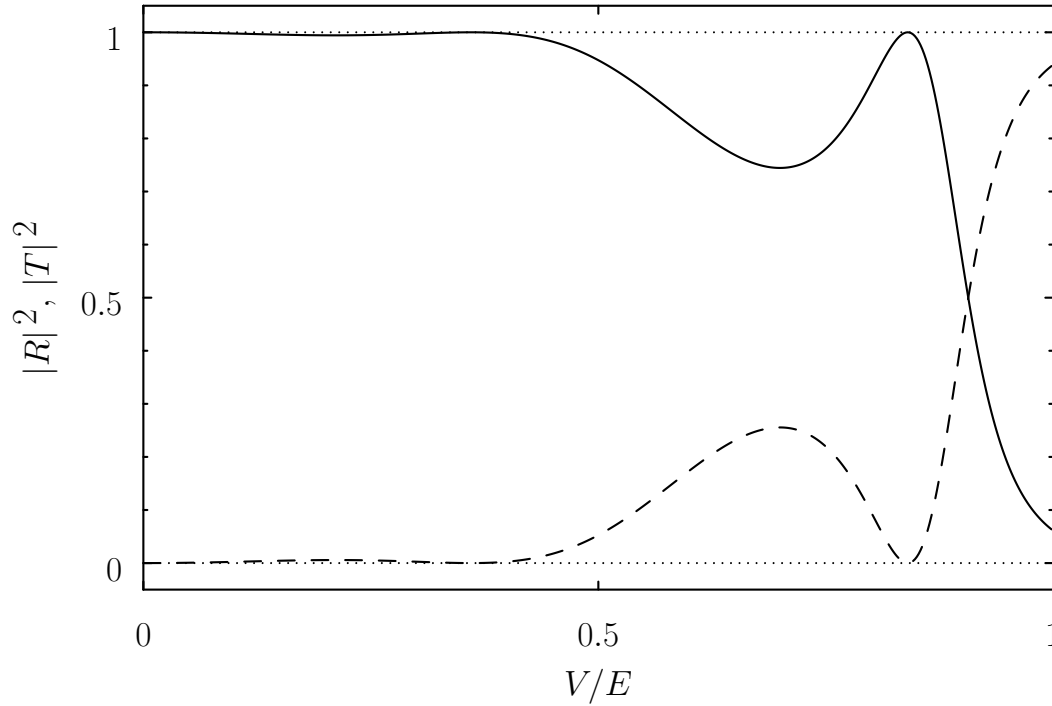


Figure 4.7: Transmission (solid curve) and reflection (dashed curve) probabilities for a square potential barrier of width $a = 1.25 \lambda$, where λ is the free-space de Broglie wavelength, as a function of the ratio of the height of the barrier, V , to the energy, E , of the incident particle.

where

$$q = \sqrt{\frac{2m(V-E)}{\hbar^2}}. \quad (4.122)$$

The general solution to Equation (4.121) takes the form

$$\psi(x) = A e^{qx} + B e^{-qx}. \quad (4.123)$$

Continuity of ψ and $d\psi/dx$ at the left edge of the barrier (i.e., $x = 0$) yields

$$1 + R = A + B, \quad (4.124)$$

$$ik(1 - R) = q(A - B). \quad (4.125)$$

Likewise, continuity of ψ and $d\psi/dx$ at the right edge of the barrier (i.e., $x = a$) gives

$$A e^{qa} + B e^{-qa} = T e^{ika}, \quad (4.126)$$

$$q(A e^{qa} - B e^{-qa}) = ikT e^{ika}. \quad (4.127)$$

After considerable algebra, the preceding four equations yield

$$|T|^2 = 1 - |R|^2 = \frac{4k^2 q^2}{4k^2 q^2 + (k^2 + q^2)^2 \sinh^2(qa)}. \quad (4.128)$$

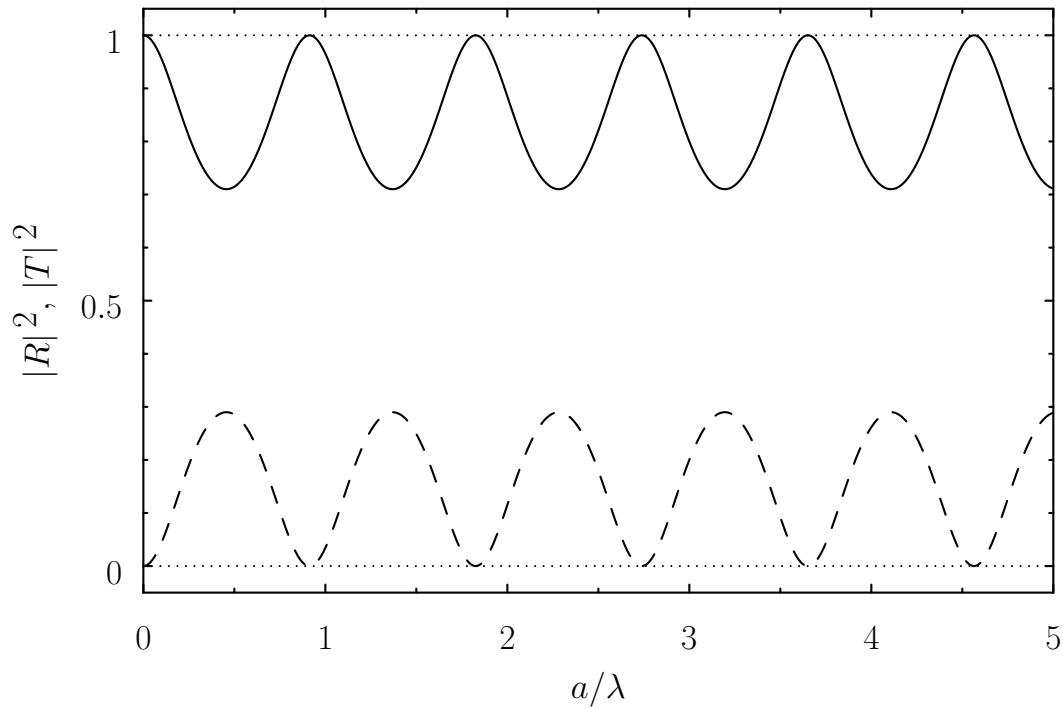


Figure 4.8: Transmission (solid curve) and reflection (dashed curve) probabilities for a particle of energy E incident on a square potential barrier of height $V = 0.75 E$ as a function of the ratio of the width of the barrier, a , to the free-space de Broglie wavelength, λ .

The fact that $|R|^2 + |T|^2 = 1$ again ensures that the probabilities of reflection and transmission sum to unity, as must be the case, because reflection and transmission are the only possible outcomes for a particle incident on the barrier.

The reflection and transmission probabilities obtained from Equation (4.128) are plotted in Figures 4.9 and 4.10. It can be seen, from these two figures, that the classical result, $|R|^2 = 1$ and $|T|^2 = 0$, is obtained for relatively thin barriers (i.e., $qa \sim 1$) in the limit where the height of the barrier is relatively large (i.e., $V \gg E$). However, if V is of order E then there is a substantial probability that the incident particle will be transmitted by the barrier. According to classical physics, transmission is impossible when $V > E$.

It can also be seen, from Figure 4.10, that the transmission probability decays exponentially as the width of the barrier increases. Nevertheless, even for very wide barriers (i.e., $qa \gg 1$), there is a small but finite probability that a particle incident on the barrier will be transmitted. This phenomenon, which is inexplicable within the context of classical physics, is called *tunneling*. For the case of a very high barrier, such that $V \gg E$, the tunneling probability reduces to

$$|T|^2 \simeq \frac{4E}{V} e^{-2a/\lambda}, \quad (4.129)$$

where $\lambda = (\hbar^2/2mV)^{1/2}$ is the de Broglie wavelength inside the barrier. Here, it is assumed that $a \gg \lambda$. Thus, even in the limit that the barrier is very high, there is an exponentially small, but

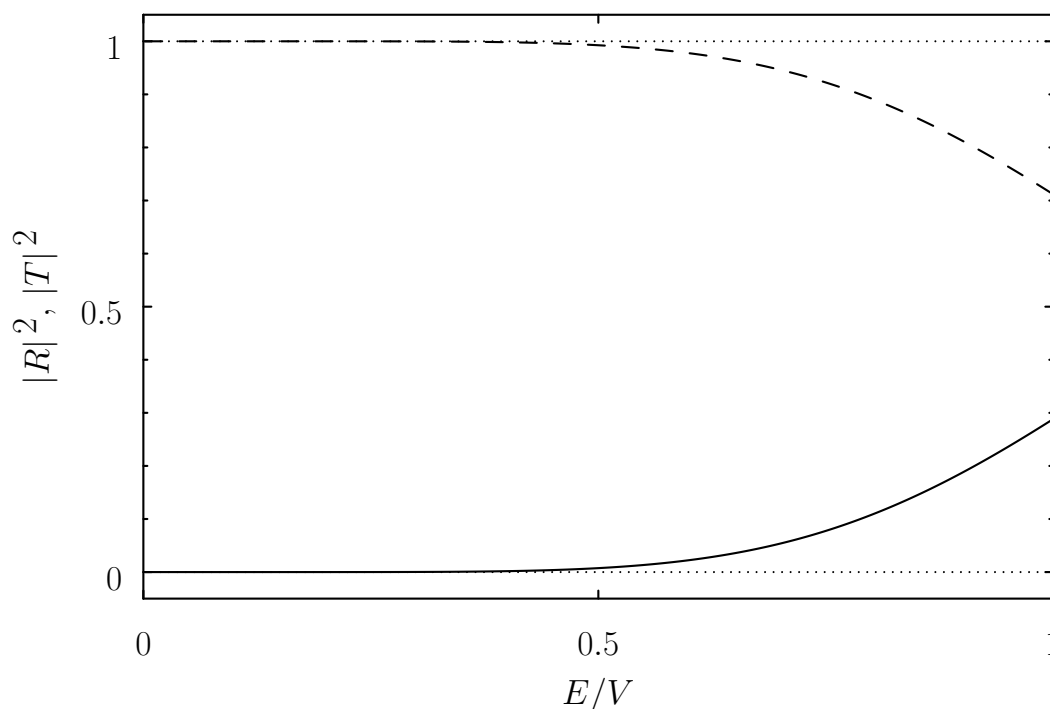


Figure 4.9: Transmission (solid curve) and reflection (dashed curve) probabilities for a square potential barrier of width $a = 0.5 \lambda$, where λ is the free-space de Broglie wavelength, as a function of the ratio of the energy, E , of the incoming particle to the height, V , of the barrier.

nevertheless non-zero, tunneling probability. Quantum mechanical tunneling plays an important role in the physics of electron field emission and α -decay. (See Sections 4.3.5 and 4.3.6.)

4.3.4 WKB Approximation

Consider a particle of mass m and energy $E > 0$ moving through some slowly-varying potential, $U(x)$. The particle's wavefunction satisfies

$$\frac{d^2\psi(x)}{dx^2} = -k^2(x)\psi(x), \quad (4.130)$$

where

$$k^2(x) = \frac{2m[E - U(x)]}{\hbar^2}. \quad (4.131)$$

Let us try a solution to Equation (4.130) of the form

$$\psi(x) = \psi_0 \exp\left(\int_0^x i k(x') dx'\right), \quad (4.132)$$

where ψ_0 is a complex constant. Note that this solution represents a particle propagating in the positive x -direction [because the full wavefunction is multiplied by $\exp(-i\omega t)$, where $\omega = E/\hbar >$

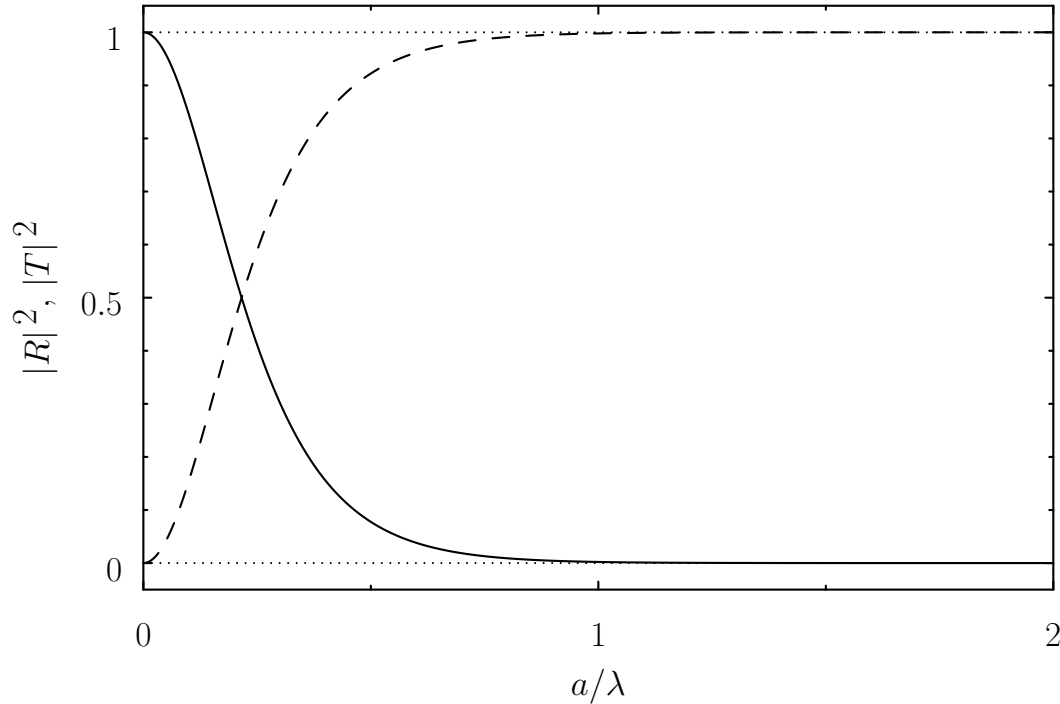


Figure 4.10: Transmission (solid curve) and reflection (dashed curve) probabilities for a particle of energy E incident on a square potential barrier of height $V = (4/3)E$ as a function of the ratio of the width of the barrier, a , to the free-space de Broglie wavelength, λ .

0] with the continuously-varying wavenumber $k(x)$. It follows that

$$\frac{d\psi(x)}{dx} = ik(x)\psi(x), \quad (4.133)$$

and

$$\frac{d^2\psi(x)}{dx^2} = ik'(x)\psi(x) - k^2(x)\psi(x), \quad (4.134)$$

where $k' \equiv dk/dx$. A comparison of Equations (4.130) and (4.134) reveals that Equation (4.132) represents an approximate solution to Equation (4.130) provided that the first term on the right-hand side of Equation (4.134) is negligible compared to the second. This yields the validity criterion $|k'| \ll k^2$, or

$$\frac{k}{|k'|} \gg k^{-1}. \quad (4.135)$$

In other words, the variation lengthscale of $k(x)$, which is approximately the same as the variation lengthscale of $U(x)$, must be much greater than the particle's de Broglie wavelength (which is of order k^{-1}). Let us suppose that this is the case. Incidentally, the approximation involved in dropping the first term on the right-hand side of Equation (4.134) is generally known as the *WKB approximation*, after G. Wentzel, H.A. Kramers, and L. Brillouin. Similarly, Equation (4.132) is termed a WKB solution.

According to the WKB solution, (4.132), the probability density remains constant; that is,

$$|\psi(x)|^2 = |\psi_0|^2; \quad (4.136)$$

as long as the particle moves through a region in which $E > U(x)$, and $k(x)$ is consequently real (i.e., an allowed region according to classical physics). Suppose, however, that the particle encounters a potential barrier (i.e., a region from which the particle is excluded according to classical physics). By definition, $E < U(x)$ inside such a barrier, and $k(x)$ is consequently imaginary. Let the barrier extend from $x = x_1$ to x_2 , where $0 < x_1 < x_2$. The WKB solution inside the barrier is written

$$\psi(x) = \psi_1 \exp\left(-\int_{x_1}^x |k(x')| dx'\right), \quad (4.137)$$

where

$$\psi_1 = \psi_0 \exp\left(\int_0^{x_1} i k(x') dx'\right). \quad (4.138)$$

Here, we have neglected the unphysical exponentially-growing solution.

According to the WKB solution, (4.137), the probability density decays exponentially inside the barrier; that is,

$$|\psi(x)|^2 = |\psi_1|^2 \exp\left(-2 \int_{x_1}^x |k(x')| dx'\right), \quad (4.139)$$

where $|\psi_1|^2$ is the probability density at the left-hand side of the barrier (i.e., $x = x_1$). It follows that the probability density at the right-hand side of the barrier (i.e., $x = x_2$) is

$$|\psi_2|^2 = |\psi_1|^2 \exp\left(-2 \int_{x_1}^{x_2} |k(x')| dx'\right). \quad (4.140)$$

Note that $|\psi_2|^2 < |\psi_1|^2$. Of course, in the region to the right of the barrier (i.e., $x > x_2$), the probability density takes the constant value $|\psi_2|^2$.

We can interpret the ratio of the probability densities to the right and to the left of the potential barrier as the probability, $|T|^2$, that a particle incident from the left will tunnel through the barrier and emerge on the other side; that is,

$$|T|^2 = \frac{|\psi_2|^2}{|\psi_1|^2} = \exp\left(-2 \int_{x_1}^{x_2} |k(x')| dx'\right). \quad (4.141)$$

It is easily demonstrated that the probability of a particle incident from the right tunneling through the barrier is the same.

Note that the criterion (4.135) for the validity of the WKB approximation implies that the previous transmission probability is very small. Hence, the WKB approximation only applies to situations in which there is very little chance of a particle tunneling through the potential barrier in question. Unfortunately, the validity criterion (4.135) breaks down completely at the edges of the barrier (i.e., at $x = x_1$ and x_2), because $k(x) = 0$ at these points. However, it can be demonstrated that the contribution of those regions, around $x = x_1$ and x_2 , in which the WKB approximation breaks down to the integral in Equation (4.141) is fairly negligible. Hence, the previous expression for the tunneling probability is a reasonable approximation provided that the incident particle's de Broglie wavelength is much smaller than the spatial extent of the potential barrier.

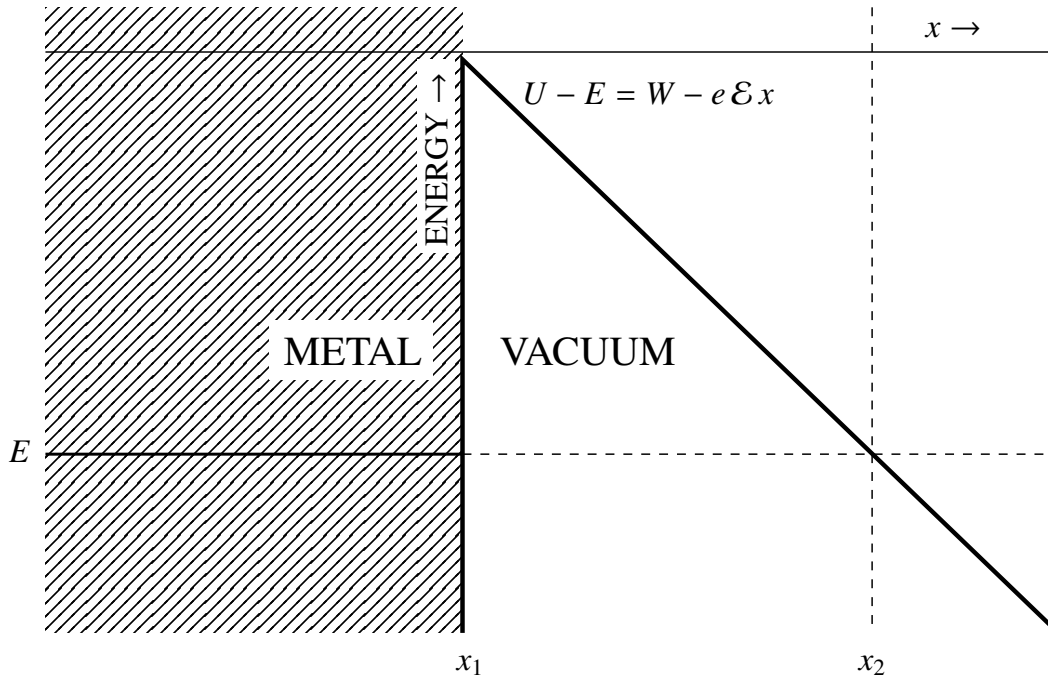


Figure 4.11: The potential barrier for an electron in a metal surface subject to an external electric field.

4.3.5 Cold Emission

Suppose that an unheated metal surface is subject to a large uniform external electric field, of strength \mathcal{E} , which is directed such that it accelerates electrons away from the surface. We have already seen (in Section 4.1.2) that electrons just below the surface of a metal can be regarded as being in a potential well of depth W , where W is called the workfunction of the surface. Adopting a simple one-dimensional treatment of the problem, let the metal lie at $x < 0$, and the surface at $x = 0$. Now, the applied electric field is shielded from the interior of the metal. (See Section 2.1.12.) Hence, the energy, E , say, of an electron just below the surface is unaffected by the field. In the absence of the electric field, the potential barrier just above the surface is simply $U(x) - E = W$. The electric field modifies this to $U(x) - E = W - e\mathcal{E}x$, where e is the magnitude of the electron charge. The potential barrier is sketched in Figure 4.11.

It can be seen, from Figure 4.11, that an electron just below the surface of the metal is confined by a triangular potential barrier that extends from $x = x_1$ to x_2 , where $x_1 = 0$ and $x_2 = W/(e\mathcal{E})$. Making use of the WKB approximation (see Section 4.3.4), the probability of such an electron tunneling through the barrier, and consequently being emitted from the surface, is

$$|T|^2 = \exp\left(-\frac{2\sqrt{2}m_e}{\hbar} \int_{x_1}^{x_2} \sqrt{U(x) - E} dx\right), \quad (4.142)$$

or

$$|T|^2 = \exp\left(-\frac{2\sqrt{2}m_e}{\hbar} \int_0^{W/e\mathcal{E}} \sqrt{W - e\mathcal{E}x} dx\right), \quad (4.143)$$

where m_e is the electron mass. This reduces to

$$|T|^2 = \exp\left(-2\sqrt{2} \frac{m_e^{1/2} W^{3/2}}{\hbar e \mathcal{E}} \int_0^1 \sqrt{1-y} dy\right), \quad (4.144)$$

or

$$|T|^2 = \exp\left(-\frac{4\sqrt{2}}{3} \frac{m_e^{1/2} W^{3/2}}{\hbar e \mathcal{E}}\right). \quad (4.145)$$

The previous result is known as the *Fowler–Nordheim formula*, after Ralph Fowler and Lothar Nordheim who derived it in 1928. Note that the probability of emission increases exponentially as the electric field-strength above the surface of the metal increases.

The cold emission of electrons from a metal surface is the basis of an important device known as a *scanning tunneling microscope*, or an STM. An STM consists of a very sharp conducting probe that is scanned over the surface of a metal (or any other solid conducting medium). A large voltage difference is applied between the probe and the surface. Now, the surface electric field-strength immediately below the probe tip is proportional to the applied potential difference, and inversely proportional to the spacing between the tip and the surface. Electrons tunneling between the surface and the probe tip give rise to a weak electric current. The magnitude of this current is proportional to the tunneling probability, (4.145). It follows that the current is an extremely sensitive function of the surface electric field-strength, and, hence, of the spacing between the tip and the surface (assuming that the potential difference is held constant). An STM can, thus, be used to construct a very accurate contour map of the surface under investigation. In fact, STMs are capable of achieving sufficient resolution to image individual atoms.

4.3.6 Alpha Decay

Many types of heavy atomic nuclei spontaneously decay to produce daughter nuclei via the emission of α -particles (i.e., helium nuclei) of some characteristic energy. This process is known as α -decay. Let us investigate the α -decay of a particular type of atomic nucleus of radius R , charge-number Z , and mass-number A . Such a nucleus thus decays to produce a daughter nucleus of charge-number $Z_1 = Z - 2$ and mass-number $A_1 = A - 4$, and an α -particle of charge-number $Z_2 = 2$ and mass-number $A_2 = 4$. Let the characteristic energy of the α -particle be E . Incidentally, nuclear radii are found to satisfy the empirical formula

$$R = 1.5 \times 10^{-15} A^{1/3} \text{ m} = 2.0 \times 10^{-15} Z_1^{1/3} \text{ m} \quad (4.146)$$

for $Z \gg 1$.

In 1928, George Gamov proposed a very successful theory of α -decay, according to which the α -particle moves freely inside the nucleus, and is emitted after tunneling through the potential barrier between itself and the daughter nucleus. In other words, the α -particle, whose energy is E , is trapped in a potential well of radius R by the potential barrier

$$U(r) = \frac{Z_1 Z_2 e^2}{4\pi \epsilon_0 r} \quad (4.147)$$

for $r > R$. (See Section 2.1.4.) Here, e is the magnitude of the electron charge.

Making use of the WKB approximation (and neglecting the fact that r is a radial, rather than a Cartesian, coordinate), the probability of the α -particle tunneling through the barrier is

$$|T|^2 = \exp \left(-\frac{2\sqrt{2}m}{\hbar} \int_{r_1}^{r_2} \sqrt{U(r) - E} dr \right), \quad (4.148)$$

where $r_1 = R$ and $r_2 = Z_1 Z_2 e^2 / (4\pi \epsilon_0 E)$. Here, $m = 4m_p$ is the α -particle mass, and m_p is the proton mass. The previous expression reduces to

$$|T|^2 = \exp \left[-2\sqrt{2}\beta \int_1^{E_c/E} \left(\frac{1}{y} - \frac{E}{E_c} \right)^{1/2} dy \right], \quad (4.149)$$

where

$$\beta = \left(\frac{Z_1 Z_2 e^2 m R}{4\pi \epsilon_0 \hbar^2} \right)^{1/2} = 0.74 Z_1^{2/3} \quad (4.150)$$

is a dimensionless constant, and

$$E_c = \frac{Z_1 Z_2 e^2}{4\pi \epsilon_0 R} = 1.44 Z_1^{2/3} \text{ MeV} \quad (4.151)$$

is the characteristic energy the α -particle would need in order to escape from the nucleus without tunneling. Of course, $E \ll E_c$. It is easily demonstrated that

$$\int_1^{1/\epsilon} \left(\frac{1}{y} - \epsilon \right)^{1/2} dy \simeq \frac{\pi}{2\sqrt{\epsilon}} - 2 \quad (4.152)$$

when $\epsilon \ll 1$. Hence,

$$|T|^2 \simeq \exp \left[-2\sqrt{2}\beta \left(\frac{\pi}{2} \sqrt{\frac{E_c}{E}} - 2 \right) \right]. \quad (4.153)$$

Now, the α -particle moves inside the nucleus at the characteristic velocity $v = \sqrt{2E/m}$. It follows that the particle bounces backward and forward within the nucleus at the frequency $\nu \simeq v/R$, giving

$$\nu \simeq 2 \times 10^{28} \text{ yr}^{-1} \quad (4.154)$$

for a 1 MeV α -particle trapped inside a typical heavy nucleus of radius 10^{-14} m. Thus, the α -particle effectively attempts to tunnel through the potential barrier ν times a second. If each of these attempts has a probability $|T|^2$ of succeeding then the probability of decay per unit time is $\nu|T|^2$. Hence, if there are $N(t) \gg 1$ intact nuclei at time t then there are only $N + dN$ at time $t + dt$, where

$$dN = -N \nu |T|^2 dt. \quad (4.155)$$

This expression can be integrated to give

$$N(t) = N(0) \exp(-\nu |T|^2 t). \quad (4.156)$$

The *half-life*, τ , is defined as the time which must elapse in order for half of the nuclei originally present to decay. It follows from the previous formula that

$$\tau = \frac{\ln 2}{\nu |T|^2}. \quad (4.157)$$

Note that the half-life is independent of $N(0)$.

Finally, making use of the previous results, we obtain

$$\log_{10}[\tau(\text{yr})] = -C_1 - C_2 Z_1^{2/3} + C_3 \frac{Z_1}{\sqrt{E(\text{MeV})}}, \quad (4.158)$$

where

$$C_1 = 28.5, \quad (4.159)$$

$$C_2 = 1.83, \quad (4.160)$$

$$C_3 = 1.73. \quad (4.161)$$

The half-life, τ , the daughter charge-number, $Z_1 = Z - 2$, and the α -particle energy, E , for atomic nuclei that undergo α -decay are indeed found to satisfy a relationship of the form (4.158). See Figure 4.12. The best fit to the data shown in the figure is obtained using

$$C_1 = 28.9, \quad (4.162)$$

$$C_2 = 1.60, \quad (4.163)$$

$$C_3 = 1.61. \quad (4.164)$$

It can be seen that these values are remarkably similar to those calculated previously.

4.3.7 Simple Harmonic Oscillator

Consider the motion of a particle of mass m in the simple harmonic oscillator potential

$$U(x) = \frac{1}{2} \kappa x^2, \quad (4.165)$$

where $\kappa > 0$ is the so-called force constant of the oscillator. According to classical physics, a particle trapped in this potential executes simple harmonic motion at the angular frequency $\omega = \sqrt{\kappa/m}$. (See Section 1.3.6.) The time-independent Schrödinger equation for a particle of mass m and energy E moving in a simple harmonic potential becomes

$$\frac{d^2\psi}{dx^2} = \frac{2m}{\hbar^2} \left(\frac{1}{2} \kappa x^2 - E \right) \psi. \quad (4.166)$$

[See Equation (4.71).] Let

$$y = \sqrt{\frac{m\omega}{\hbar}} x, \quad (4.167)$$

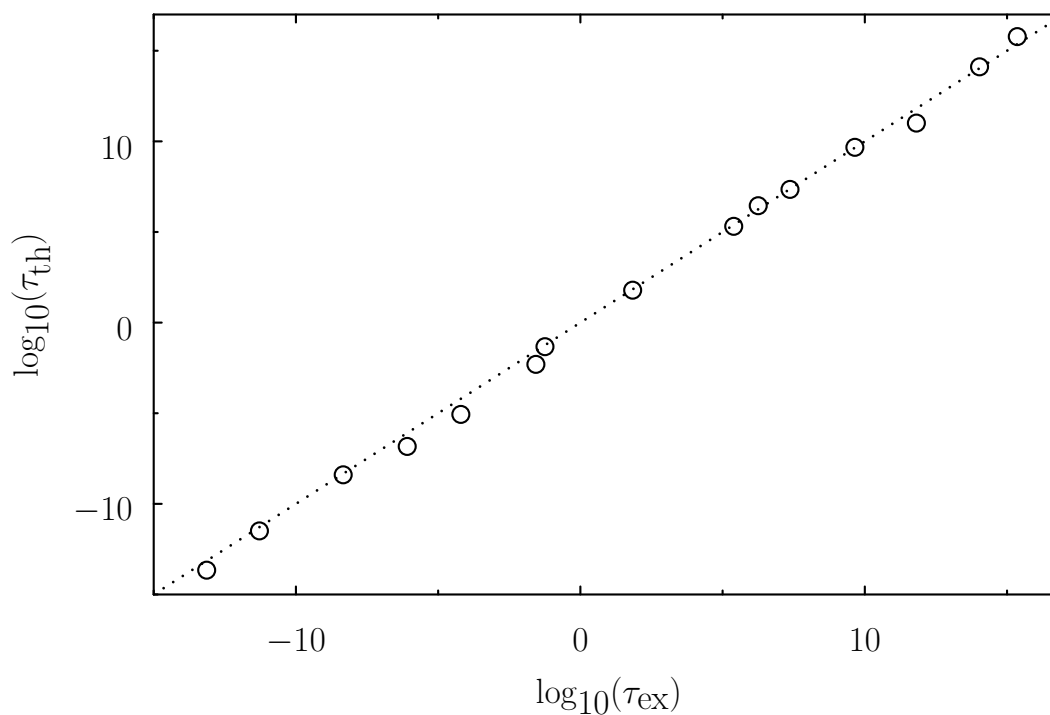


Figure 4.12: The experimentally determined half-life, τ_{ex} , of various atomic nuclei that decay via α -emission versus the best-fit theoretical half-life $\log_{10}(\tau_{\text{th}}) = -28.9 - 1.60 Z_1^{2/3} + 1.61 Z_1 / \sqrt{E}$. Both half-lives are measured in years. Here, $Z_1 = Z - 2$, where Z is the charge-number of the nucleus, and E the characteristic energy of the emitted α -particle in MeV. In order of increasing half-life, the points correspond to the following nuclei: Rn 215, Po 214, Po 216, Po 197, Fm 250, Ac 225, U 230, U 232, U 234, Gd 150, U 236, U 238, Pt 190, Gd 152, Nd 144. (Data obtained from International Atomic Energy Agency, Nuclear Data Center.)

and

$$\epsilon = \frac{2E}{\hbar\omega}. \quad (4.168)$$

Equation (4.166) reduces to

$$\frac{d^2\psi}{dy^2} - (y^2 - \epsilon)\psi = 0. \quad (4.169)$$

We need to find solutions to the previous equation that are bounded at infinity; that is, solutions that satisfy the boundary condition $\psi \rightarrow 0$ as $|y| \rightarrow \infty$.

Consider the behavior of the solution to Equation (4.169) in the limit $|y| \gg 1$. As is easily seen, in this limit the equation simplifies somewhat to give

$$\frac{d^2\psi}{dy^2} - y^2\psi \simeq 0. \quad (4.170)$$

The approximate solutions to the previous equation are

$$\psi(y) \simeq A(y) e^{\pm y^2/2}, \quad (4.171)$$

where $A(y)$ is a relatively slowly varying function of y . Clearly, if $\psi(y)$ is to remain bounded as $|y| \rightarrow \infty$ then we must choose the exponentially decaying solution. This suggests that we should write

$$\psi(y) = h(y) e^{-y^2/2}, \quad (4.172)$$

where we would expect $h(y)$ to be an algebraic, rather than an exponential, function of y .

Substituting Equation (4.172) into Equation (4.169), we obtain

$$\frac{d^2h}{dy^2} - 2y \frac{dh}{dy} + (\epsilon - 1)h = 0. \quad (4.173)$$

Let us attempt a power-law solution of the form

$$h(y) = \sum_{i=0, \infty} c_i y^i. \quad (4.174)$$

Inserting this test solution into Equation (4.173), and equating the coefficients of y^i , we obtain the recursion relation

$$c_{i+2} = \frac{(2i - \epsilon + 1)}{(i+1)(i+2)} c_i. \quad (4.175)$$

Consider the behavior of $h(y)$ in the limit $|y| \rightarrow \infty$. The previous recursion relation simplifies to

$$c_{i+2} \simeq \frac{2}{i} c_i. \quad (4.176)$$

Hence, at large $|y|$, when the higher powers of y dominate, we have

$$h(y) \sim C \sum_j \frac{y^{2j}}{j!} \sim C e^{y^2}. \quad (4.177)$$

It follows that $\psi(y) = h(y) \exp(-y^2/2)$ varies as $\exp(y^2/2)$ as $|y| \rightarrow \infty$. This behavior is unacceptable, because it does not satisfy the boundary condition $\psi \rightarrow 0$ as $|y| \rightarrow \infty$. The only way in which we can prevent ψ from blowing up as $|y| \rightarrow \infty$ is to demand that the power series (4.174) terminate at some finite value of i . This implies, from the recursion relation (4.175), that

$$\epsilon = 2n + 1, \quad (4.178)$$

where n is a non-negative integer. Note that the number of terms in the power series (4.174) is $n + 1$. Finally, using Equation (4.168), we obtain

$$E = \left(\frac{1}{2} + n \right) \hbar \omega, \quad (4.179)$$

for $n = 0, 1, 2, \dots$.

Hence, we conclude that a particle moving in a harmonic potential has quantized energy levels that are equally spaced. The spacing between successive energy levels is $\hbar \omega$, where ω is the classical oscillation frequency. Furthermore, the lowest energy state ($n = 0$) possesses the finite energy $(1/2) \hbar \omega$. This is another example of zero-point energy. (See Section 4.3.1.) It is easily demonstrated that the (normalized) wavefunction of the lowest energy state takes the form

$$\psi_0(x) = \frac{e^{-x^2/2d^2}}{\pi^{1/4} \sqrt{d}}, \quad (4.180)$$

where $d = \sqrt{\hbar/m\omega}$.

4.4 Three-Dimensional Wave Mechanics

4.4.1 Three-Dimensional Wave Mechanics

Up to now, we have only discussed wave mechanics for a particle moving in one dimension. However, the generalization to a particle moving in three dimensions is fairly straightforward. A massive particle moving in three dimensions has a complex wavefunction of the form [cf., Equation (4.10)]

$$\psi(x, y, z, t) = \psi_0 e^{i(\mathbf{k}\cdot\mathbf{r} - \omega t)}, \quad (4.181)$$

where ψ_0 is a complex constant, and $\mathbf{r} = (x, y, z)$. Here, the wavevector, \mathbf{k} , and the angular frequency, ω , are related to the particle momentum, \mathbf{p} , and energy, E , according to [cf., Equation (4.9)]

$$\mathbf{p} = \hbar \mathbf{k}, \quad (4.182)$$

and [cf., Equation (4.8)]

$$E = \hbar \omega, \quad (4.183)$$

respectively. Generalizing the analysis of Section 4.2.2, the three-dimensional version of Schrödinger's equation is [cf., Equation (4.22)]

$$i \hbar \frac{\partial \psi}{\partial t} = -\frac{\hbar^2}{2m} \nabla^2 \psi + U(\mathbf{r}) \psi, \quad (4.184)$$

where the differential operator

$$\nabla^2 \equiv \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} \quad (4.185)$$

is known as the *Laplacian*. (See Section A.21.) The interpretation of a three-dimensional wavefunction is that the probability of simultaneously finding the particle between x and $x + dx$, between y and $y + dy$, and between z and $z + dz$, at time t is [cf., Equation (4.25)]

$$P(x, y, z, t) = |\psi(x, y, z, t)|^2 dx dy dz. \quad (4.186)$$

Moreover, the normalization condition for the wavefunction becomes [cf., Equation (4.27)]

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |\psi(x, y, z, t)|^2 dx dy dz = 1. \quad (4.187)$$

It can be demonstrated that Schrödinger's equation, (4.184), preserves the normalization condition, (4.187), of a localized wavefunction. Heisenberg's uncertainty principle generalizes to [cf., Equation (4.65)]

$$\Delta x \Delta p_x \gtrsim \frac{\hbar}{2}, \quad (4.188)$$

$$\Delta y \Delta p_y \gtrsim \frac{\hbar}{2}, \quad (4.189)$$

$$\Delta z \Delta p_z \gtrsim \frac{\hbar}{2}. \quad (4.190)$$

Finally, a stationary state of energy E is written [cf., Equation (4.69)]

$$\psi(x, y, z, t) = \psi(x, y, z) e^{-i(E/\hbar)t}, \quad (4.191)$$

where the stationary wavefunction, $\psi(x, y, z)$, satisfies [cf., Equation (4.71)]

$$-\frac{\hbar^2}{2m} \nabla^2 \psi + U(\mathbf{r}) \psi = E \psi. \quad (4.192)$$

4.4.2 Particle in Box

As an example of a three-dimensional problem in wave mechanics, consider a particle trapped in a square potential well of infinite depth, such that

$$U(x, y, z) = \begin{cases} 0 & 0 \leq x \leq a, 0 \leq y \leq a, 0 \leq z \leq a \\ \infty & \text{otherwise} \end{cases}. \quad (4.193)$$

Within the well, the stationary wavefunction, $\psi(x, y, z)$, satisfies

$$-\frac{\hbar^2}{2m} \nabla^2 \psi = E \psi, \quad (4.194)$$

subject to the boundary conditions

$$\psi(0, y, z) = \psi(x, 0, z) = \psi(x, y, 0) = 0, \quad (4.195)$$

and

$$\psi(a, y, z) = \psi(x, a, z) = \psi(x, y, a) = 0, \quad (4.196)$$

because $\psi = 0$ outside the well. Let us try a separable wavefunction of the form

$$\psi(x, y, z) = \psi_0 \sin(k_x x) \sin(k_y y) \sin(k_z z). \quad (4.197)$$

This expression automatically satisfies the boundary conditions (4.195). The remaining boundary conditions, (4.196), are satisfied provided

$$k_x = n_x \frac{\pi}{a}, \quad (4.198)$$

$$k_y = n_y \frac{\pi}{a}, \quad (4.199)$$

$$k_z = n_z \frac{\pi}{a}, \quad (4.200)$$

where n_x , n_y , and n_z are (independent) positive integers. (Note that a negative value of n_x does not give rise to a physical state that is distinct from the corresponding positive value, et cetera.) Substitution of the wavefunction (4.197) into Equation (4.194) yields

$$E = \frac{\hbar^2 k^2}{2m} = \frac{\hbar^2}{2m} (k_x^2 + k_y^2 + k_z^2). \quad (4.201)$$

Thus, it follows from Equations (4.198)–(4.200) that the particle energy is quantized, and that the allowed *energy levels* are

$$E_{l_x, l_y, l_z} = \frac{\hbar^2}{2m a^2} (n_x^2 + n_y^2 + n_z^2). \quad (4.202)$$

The properly normalized [see Equation (4.187)] stationary wavefunctions corresponding to these energy levels are

$$\psi_{l_x, l_y, l_z}(x, y, z) = \left(\frac{2}{a}\right)^{3/2} \sin\left(n_x \pi \frac{x}{a}\right) \sin\left(n_y \pi \frac{y}{a}\right) \sin\left(n_z \pi \frac{z}{a}\right). \quad (4.203)$$

As is the case for a particle trapped in a one-dimensional potential well, the lowest energy level for a particle trapped in a three-dimensional well is not zero, but rather

$$E_{1,1,1} = 3 E_1. \quad (4.204)$$

Here,

$$E_1 = \frac{\hbar^2}{2m a^2}. \quad (4.205)$$

is the *ground state* (i.e., the lowest energy state) energy in the one-dimensional case. It follows from Equation (4.202) that distinct permutations of n_x , n_y , and n_z that do not alter the value of $n^2 = n_x^2 + n_y^2 + n_z^2$ also do not alter the energy. In other words, in three dimensions, it is possible for distinct wavefunctions to be associated with the same energy level. In this situation, the energy level is said to be *degenerate*. The ground-state energy level, $3 E_1$, is non-degenerate, because the only combination of (n_x, n_y, n_z) that gives this energy is $(1, 1, 1)$. However, the next highest energy level, $6 E_1$, is degenerate, because it is obtained when (n_x, n_y, n_z) take the values $(2, 1, 1)$, or $(1, 2, 1)$, or $(1, 1, 2)$. In fact, a non-degenerate energy level corresponds to a case where the three quantum numbers (i.e., n_x , n_y , and n_z) all have the same value, whereas a threefold degenerate energy level corresponds to a case where only two of the quantum numbers have the same value, and, finally, a sixfold degenerate energy level corresponds to a case where the quantum numbers are all different.

4.4.3 Degenerate Electron Gas

Consider N electrons trapped in a cubic box of dimension a . Let us treat the electrons as essentially non-interacting particles. The total energy of a system consisting of many non-interacting particles is simply the sum of the single-particle energies of the individual particles. Furthermore, because the electrons are indistinguishable fermions (i.e., half-integer spin particle), they are subject to the so-called *Pauli exclusion principle*. The exclusion principle states that no two electrons in our system can occupy the same single-particle energy level. Now, from Section 4.4.2, the single-particle energy levels for a particle in a box are characterized by the three quantum numbers, n_x , n_y , and n_z . Thus, we conclude that no two electrons in our system can have the same set of values of n_x , n_y , and n_z . It turns out that this is not quite true, because electrons possess an intrinsic angular momentum called *spin*. The spin states of an electron are governed by an additional quantum number that can take one of two different values. Hence, when spin is taken into account, we conclude that a maximum of two electrons (with different spin quantum numbers) can occupy a single-particle energy level corresponding to a particular set of values of n_x , n_y , and n_z . It is clear, from Equation (4.202), that the associated particle energy is proportional to $n^2 = n_x^2 + n_y^2 + n_z^2$.

Suppose that our electrons are cold; that is, they have comparatively little thermal energy. In this case, we would expect them to fill the lowest single-particle energy levels available to them. We can imagine the single-particle energy levels as existing in a sort of three-dimensional quantum number space whose Cartesian coordinates are n_x , n_y , and n_z . Thus, the energy levels are uniformly distributed in this space on a cubic lattice. Moreover, the distance between nearest-neighbor energy levels is unity. This implies that the number of energy levels per unit volume is also unity. Finally, the energy of a given energy level is proportional to its distance, $n^2 = n_x^2 + n_y^2 + n_z^2$, from the origin.

Because we expect cold electrons to occupy the lowest energy levels available to them, but only two electrons can occupy a given energy level, it follows that if the number of electrons, N_e , is very large then the filled energy levels will be approximately distributed in a sphere centered on the origin of quantum number space. The number of energy levels contained in a sphere of radius n is approximately equal to the volume of the sphere, because the number of energy levels per unit volume is unity. It turns out that this is not quite correct, because we have forgotten that the quantum numbers n_x , n_y , and n_z can only take positive values. Hence, the filled energy levels

actually only occupy one octant of a sphere. The radius, n_F , of the octant of filled energy levels in quantum number space can be calculated by equating the number of energy levels it contains to the number of electrons, N_e . Thus, we can write

$$N_e = 2 \times \frac{1}{8} \times \frac{4\pi}{3} n_F^3. \quad (4.206)$$

Here, the factor 2 is to take into account the two spin states of an electron, and the factor 1/8 is to take account of the fact that n_x , n_y , and n_z can only take positive values. Thus,

$$n_F = \left(\frac{3 N_e}{\pi} \right)^{1/3}. \quad (4.207)$$

According to Equation (4.202), the energy of the most energetic electrons—which is known as the *Fermi energy*—is given by

$$E_F = \frac{n_F^2 \pi^2 \hbar^2}{2 m_e a^2} = \frac{\pi^2 \hbar^2}{2 m_e a^2} \left(\frac{3 N_e}{\pi} \right)^{2/3}, \quad (4.208)$$

where m_e is the electron mass. This expression can also be written as

$$E_F = \frac{\pi^2 \hbar^2}{2 m_e} \left(\frac{3 n_e}{\pi} \right)^{2/3}, \quad (4.209)$$

where $n_e = N_e/a^3$ is the number of electrons per unit volume (in real space). Note that the Fermi energy only depends on the number density of the confined electrons.

The mean energy of the electrons is given by

$$\langle E \rangle = E_F \int_0^{n_F} n^2 4\pi n^2 dn \bigg/ \frac{4}{3} \pi n_F^3 = \frac{3}{5} E_F, \quad (4.210)$$

because $E \propto n^2$, and the energy levels are uniformly distributed in quantum-number space within an octant of radius n_F . According to classical physics, the mean thermal energy of the electrons is $(3/2)k_B T$, where T is the electron temperature, and k_B the Boltzmann constant. Thus, if $k_B T \ll E_F$ then our original assumption that the electrons are cold is valid. Note that, in this case, the electron energy is much larger than that predicted by classical physics; electrons in this state are termed *degenerate*. On the other hand, if $k_B T \gg E_F$ then the electrons are hot, and are essentially governed by classical physics; electrons in this state are termed *non-degenerate*.

The total energy of a degenerate electron gas is

$$E_{\text{total}} = N_e \langle E \rangle = \frac{3}{5} N_e E_F. \quad (4.211)$$

Hence, the gas pressure takes the form

$$P = -\frac{\partial E_{\text{total}}}{\partial V} = \frac{2}{5} n E_F, \quad (4.212)$$

because $E_F \propto a^{-2} = V^{-2/3}$. [See Equation (4.208).] Now, the pressure predicted by classical physics is $P = nk_B T$. Thus, a degenerate electron gas has a much higher pressure than that which would be predicted by classical physics. This is an entirely quantum mechanical effect, and is due to the fact that identical fermions cannot get significantly closer together than a de Broglie wavelength without violating the Pauli exclusion principle. Note that, according to Equation (4.209), the mean spacing between degenerate electrons is

$$d \sim n_e^{-1/3} \sim \frac{h}{\sqrt{m_e E}} \sim \frac{h}{p} \sim \lambda, \quad (4.213)$$

where λ is the de Broglie wavelength. Thus, an electron gas is non-degenerate when the mean spacing between the electrons is much greater than the de Broglie wavelength, and becomes degenerate as the mean spacing approaches the de Broglie wavelength.

It turns out that the conduction (i.e., free) electrons inside metals are highly degenerate (because the number of electrons per unit volume is very large, and $E_F \propto n_e^{2/3}$). Indeed, most metals are hard to compress as a direct consequence of the high degeneracy pressure of their conduction electrons. To be more exact, resistance to compression is usually measured in terms of a quantity known as the *bulk modulus*, which is defined

$$\kappa = -V \frac{\partial P}{\partial V} \quad (4.214)$$

Now, for a fixed number of electrons, $P \propto V^{-5/3}$ [see Equations (4.208) and (4.212)]. Hence,

$$\kappa = \frac{5}{3} P = \frac{\pi^3 \hbar^2}{9 m_e} \left(\frac{3 n_e}{\pi} \right)^{5/3}. \quad (4.215)$$

For example, the number density of free electrons in magnesium is $n_e \sim 8.6 \times 10^{28} \text{ m}^{-3}$. This leads to the following estimate for the bulk modulus; $\kappa \sim 6.4 \times 10^{10} \text{ N m}^{-2}$. The actual bulk modulus is $\kappa = 4.5 \times 10^{10} \text{ N m}^{-2}$.

4.4.4 White-Dwarf Star

A main-sequence hydrogen-burning star, such as the Sun, is maintained in equilibrium via the balance of the gravitational attraction tending to make it collapse, and the thermal pressure tending to make it expand. Of course, the thermal energy of the star is generated by nuclear reactions occurring deep inside its core. Eventually, however, the star will run out of burnable fuel, and, therefore, start to collapse, as it radiates away its remaining thermal energy. What is the ultimate fate of such a star?

A burnt-out star is basically a gas of electrons and ions. As the star collapses, its density increases, and so the mean separation between its constituent particles decreases. Eventually, the mean separation becomes of order of the de Broglie wavelength of the electrons, and the electron gas becomes degenerate. Note that the de Broglie wavelength of the ions is much smaller than that of the electrons (because the ions are much more massive), so the ion gas remains non-degenerate. Now, even at zero temperature, a degenerate electron gas exerts a substantial pressure, because

the Pauli exclusion principle prevents the mean electron separation from becoming significantly smaller than the typical de Broglie wavelength. (See Section 4.4.3.) Thus, it is possible for a burnt-out star to maintain itself against complete collapse under gravity via the degeneracy pressure of its constituent electrons. Such stars are termed *white-dwarfs*. Let us investigate the physics of white-dwarfs in more detail.

The total energy of a white-dwarf star can be written

$$\mathcal{E} = K + U, \quad (4.216)$$

where K is the kinetic energy of the degenerate electrons (the kinetic energy of the ions is negligible), and U is the gravitational potential energy. Let us assume, for the sake of simplicity, that the density of the star is uniform. In this case, the gravitational potential energy takes the form

$$U = -\frac{3}{5} \frac{G M^2}{R}, \quad (4.217)$$

where G is the gravitational constant, M is the stellar mass, and R is the stellar radius. The previous equation follows by analogy with Equation (2.90).

From the previous section, the kinetic energy of a degenerate electron gas is simply

$$K = N_e \langle E \rangle = \frac{3}{5} N E_F = \frac{3}{5} N_e \frac{\pi^2 \hbar^2}{2 m_e} \left(\frac{3 N_e}{\pi V} \right)^{2/3}, \quad (4.218)$$

where N_E is the number of electrons, V the volume of the star, and m_e the electron mass.

The interior of a white-dwarf star is composed of atoms like C^{12} and O^{16} which contain equal numbers of protons, neutrons, and electrons. Thus,

$$M = 2 N_e m_p, \quad (4.219)$$

where m_p is the proton mass.

Equations (4.216)–(4.219) can be combined to give

$$\mathcal{E} = \frac{A}{R^2} - \frac{B}{R}, \quad (4.220)$$

where

$$A = \frac{3}{20} \left(\frac{9\pi}{8} \right)^{2/3} \frac{\hbar^2}{m_e} \left(\frac{M}{m_p} \right)^{5/3}, \quad (4.221)$$

$$B = \frac{3}{5} G M^2. \quad (4.222)$$

The equilibrium radius of the star, R_* , is that which minimizes the total energy \mathcal{E} . In fact, it is easily demonstrated that

$$R_* = \frac{2A}{B}, \quad (4.223)$$

which yields

$$R_* = \frac{(9\pi)^{2/3}}{8} \frac{\hbar^2}{G m_e m_p^{5/3} M^{1/3}}. \quad (4.224)$$

The previous formula can also be written

$$\frac{R_*}{R_\odot} = 0.010 \left(\frac{M_\odot}{M} \right)^{1/3}, \quad (4.225)$$

where $R_\odot = 7 \times 10^5$ km is the solar radius, and $M_\odot = 2 \times 10^{30}$ kg the solar mass. It follows that the radius of a typical solar-mass white-dwarf is about 7000 km; that is, about the same as the radius of the Earth. The first white-dwarf to be discovered (in 1862) was the companion of Sirius. Nowadays, thousands of white-dwarfs have been observed, all with properties similar to those described previously.

Note from Equations (4.218), (4.219), and (4.225) that $\langle E \rangle \propto M^{4/3}$. In other words, the mean energy of the electrons inside a white-dwarf increases as the stellar mass increases. Hence, for a sufficiently massive white-dwarf, the electrons can become relativistic. It turns out that the degeneracy pressure for relativistic electrons only scales as R^{-1} , rather than R^{-2} , and, thus, is unable to balance the gravitational pressure (which also scales as R^{-1}). It follows that electron degeneracy pressure is only able to halt the collapse of a burnt-out star provided that the stellar mass does not exceed some critical value, known as the *Chandrasekhar limit*, because it was first derived by Subrahmanyan Chandrasekhar in 1930, which turns out to be about 1.4 times the mass of the Sun. Stars whose mass exceeds the Chandrasekhar limit inevitably collapse to produce extremely compact objects, such as neutron stars (which are held up by the degeneracy pressure of their constituent neutrons), or black holes.

Chapter 5

Thermal Physics

5.1 Probability Theory

5.1.1 Probability

Consider some physical system A . Suppose that a measurement of a given property of this system can result in a number of distinct outcomes. If we wish to determine the *probability* of obtaining a given outcome at an arbitrary time then we can take one of two approaches. First, we can observe system A at many distinct times; this approach is known as a *time average*. Second, we can observe many systems that are identical to A at an arbitrary time; this approach is known as an *ensemble average*. An ensemble average is the most convenient theoretical approach, and the one that we shall adopt in the following discussion, whereas a time average is more directly related to real experiments.

Suppose that there are N systems in our ensemble (i.e., collection of identical systems) and that N_r of these systems exhibit the outcome r . The *probability* of occurrence of outcome r is defined

$$P_r = \lim_{N \rightarrow \infty} \frac{N_r}{N}. \quad (5.1)$$

It is clear that P_r is a number that lies between 0 and 1. If $P_r = 0$ then no systems in the ensemble exhibit the outcome r , even in the limit that the number of systems tends to infinity. This is another way of saying that outcome r is impossible. If $P_r = 1$ then all systems in the ensemble exhibit the outcome r , even in the limit that the number of systems tends to infinity. This is another way of saying that outcome r is certain to occur.

Suppose that a measurement of a given property of some physical system A can lead to any one of R mutually exclusive outcomes. Let the total number of systems in the ensemble be N , and let the number of systems that exhibit the outcome r be N_r . It follows that

$$\sum_{r=1,R} N_r = N. \quad (5.2)$$

However, if we divide both sides of the previous equation by N , and then take the limit that $N \rightarrow \infty$,

then we obtain the so-called *normalization condition*,

$$\sum_{r=1,R} P_r = 1, \quad (5.3)$$

where use has been made of Equation (5.1). The normalization condition states that the sum of the probabilities of all of the possible outcomes of a measurement of a given property of system A is unity. This condition is equivalent to the self-evident proposition that a measurement of the property is bound to result in one of the possible outcomes of this measurement.

Let us determine the probability of occurrence of outcome r or outcome s when an observation is made of our system. Here, r and s are distinct outcomes. There are $N_r + N_s$ systems in our ensemble that exhibit either the outcome r or the outcome s , so

$$P_{r|s} = \lim_{N \rightarrow \infty} \frac{N_r + N_s}{N} = P_r + P_s, \quad (5.4)$$

where use has been made of Equation (5.1). In other words, the probability of observing the outcome r or the outcome s is the sum of the probabilities of occurrence of these two outcomes. For example, the probability of throwing a 1 on a six-sided die is $1/6$. Likewise, the probability of throwing a 2 is $1/6$. Hence, the probability of throwing a 1 or a 2 is $1/6 + 1/6 = 1/3$. The previous result can easily be extended to deal with more than two alternative outcomes.

Suppose that our system can exhibit two different types of outcome. Type-1 outcomes are labeled $r = 1, \dots, R$. Type-2 outcomes are labeled $s = 1, \dots, S$. Let there be N systems in our ensemble, and let N_r of them exhibit the type-1 outcome r , and let N_s of them exhibit the type-2 outcome s . The probability of outcome s is

$$P_s = \frac{N_s}{N}, \quad (5.5)$$

which implies that

$$N_s = P_s N. \quad (5.6)$$

[Here, the limit $N \rightarrow \infty$ is taken as read; see Equation (5.1).] By analogy, the number of systems that exhibit the type-1 outcome r and the type-2 outcome s is

$$N_{r \otimes s} = P_s N_r. \quad (5.7)$$

Hence, the probability of obtaining both the type-1 outcome r and the type-2 outcome s simultaneously is

$$P_{r \otimes s} = \lim_{N \rightarrow \infty} \frac{N_{r \otimes s}}{N} = P_s \lim_{N \rightarrow \infty} \frac{N_r}{N} = P_r P_s, \quad (5.8)$$

where use has been made of Equation (5.1). However, the previous result is only valid provided outcomes r and s are statistically independent of one another. In other words, obtaining the outcome r must not affect the probability of obtaining the outcome s . As an example of the previous result, consider a system consisting of two six-sided dice. The probability of throwing a 1 on either die is $1/6$. Hence, the probability of simultaneously throwing a 1 on both dice is $1/6 \times 1/6 = 1/36$. The previous result can easily be extended to deal with more than two types of outcome.

5.1.2 Binomial Probability Distribution

Consider a system that can only exhibit two possible outcomes. Let us label the outcomes P and Q , and let p and q be the respective probabilities of these outcomes. It follows from Equation (5.3) that

$$p + q = 1. \quad (5.9)$$

Consider an ensemble of N two-outcome systems like the one just discussed. Let n be the number of systems in the ensemble that exhibit outcome P , and let n' be the number of systems that exhibit outcome Q . It is evident that

$$n + n' = N. \quad (5.10)$$

Let us determine the probability, $P_N(n)$, that n systems in our ensemble exhibit outcome P . Making use of a straightforward extension of Equation (5.8), the probability that n systems in the ensemble exhibit outcome P , and that n' exhibit outcome Q , is

$$\underbrace{p p p p \cdots p}_n \underbrace{q q q q \cdots q}_{n'} = p^n q^{n'}. \quad (5.11)$$

However, a situation in which n systems in the ensemble exhibit the outcome P can be achieved in many alternative ways. Let $C_N(n)$ be the number of distinct configurations of N systems by which n of these systems exhibit outcome P . Making use of a straightforward extension of Equation (5.4), as well as Equations (5.10) and (5.11), we deduce that

$$P(n) = C_N(n) p^n q^{N-n}. \quad (5.12)$$

Consider n systems exhibiting the outcome P . The number of ways that these systems can be distributed between N systems is

$$N(N-1)(N-2)\cdots[N-(n-1)] = \frac{N!}{(N-n)!}. \quad (5.13)$$

This follows, by induction from Equation (5.8), because we can choose any one of the N systems to exhibit the first outcome P , then we can choose any one of the remaining $N-1$ systems to exhibit the second outcome P , and so on. However, some of the $N!/(N-n)!$ distributions will just be permutations of the systems exhibiting outcome P among themselves. Such permutations do not correspond to distinct distributions. Now, the number of permutations of n quantities among n places is $n!$. Hence, we deduce that

$$C_N(n) = \frac{N!}{n!(N-n)!}. \quad (5.14)$$

It follows from Equation (5.12) that

$$P_N(n) = \frac{N!}{n!(N-n)!} p^n q^{N-n}. \quad (5.15)$$

The well-known algebraic expansion of a binomial of the form $(p + q)^N$ is

$$(p + q)^N = \sum_{n=1, N} \frac{N!}{n!(N - n)!} p^n q^{N-n}. \quad (5.16)$$

For this reason, the probability distribution (5.15) is known as the *binomial probability distribution*. From Equation (5.9), $p + q = 1$, which implies that $(p + q)^N = 1$. Thus, the previous two equations yield

$$\sum_{n=0, N} P_N(n) = 1, \quad (5.17)$$

in accordance with Equation (5.3).

Suppose that outcomes P and Q represent steps to the right and steps to the left taken by a drunken man. The net number of steps to the right taken is $m = n - n' = 2n - N$, where use has been made of Equation (5.10). Thus,

$$n = \frac{N + m}{2}, \quad (5.18)$$

which implies that the probability, $P'_N(m)$, that m assumes a certain value after N steps is equal to the probability that n assumes the value $(N + m)/2$. In other words,

$$P'_N(m) = P_N\left(\frac{N + m}{2}\right). \quad (5.19)$$

Suppose, finally, that some physical system A can exhibit many possible outcomes, r, s, t , et cetera. If we are only interested in outcome r then we could label all of the other outcomes 'not r ' or \bar{r} . In this case, we have recovered a system to which the binomial probability distribution applies.

5.1.3 Mean, Variance, and Standard Deviation

Consider some physical system A . Suppose that a measurement of a particular physical property of this system, x , can result in one of R distinct outcomes. Let outcome r be associated with x taking the value x_r . Consider an ensemble of N systems that are identical to system A . Let N_r be the number of systems in the ensemble that exhibit the outcome r . The *mean value* of x is, by definition, the average of a very large number of measurements of x . In other words,

$$\langle x \rangle = \lim_{N \rightarrow \infty} \sum_{r=1, R} \frac{N_r x_r}{N}, \quad (5.20)$$

which implies that

$$\langle x \rangle = \sum_{r=1, R} P_r x_r, \quad (5.21)$$

where use has been made of Equation (5.1).

Let

$$\Delta x = x - \langle x \rangle \quad (5.22)$$

measure the deviation of an individual measurement of x from the mean value. Obviously,

$$\langle \Delta x \rangle = \langle x \rangle - \langle \langle x \rangle \rangle = \langle x \rangle - \langle x \rangle = 0. \quad (5.23)$$

In other words, the mean deviation from the mean value is zero.

Consider $\langle (\Delta x)^2 \rangle$. This quantity, which is known as the *variance* of x , is positive definite. It can only take the value 0 if all measurement of x result in the mean value. Thus, the variance of x measures the degree of scatter about the mean value. It follows that

$$\langle (\Delta x)^2 \rangle = \langle (x - \langle x \rangle)^2 \rangle = \langle x^2 - 2 \langle x \rangle x + \langle x \rangle^2 \rangle = \langle x^2 \rangle - 2 \langle x \rangle^2 + \langle x \rangle^2, \quad (5.24)$$

or

$$\langle (\Delta x)^2 \rangle = \langle x^2 \rangle - \langle x \rangle^2. \quad (5.25)$$

Finally, the quantity

$$\sigma_x = [\langle (\Delta x)^2 \rangle]^{1/2} \quad (5.26)$$

is known as the *standard deviation* of x . The standard deviation is essentially the width of the range of probable values over which x is distributed around its mean value, $\langle x \rangle$.

5.1.4 Application to Binomial Probability Distribution

Let us now apply what we have just learned about the mean, variance, and standard deviation of a general probability distribution to the specific case of the binomial probability distribution. Recall, from Section 5.1.2, that if a simple system has just two possible outcomes, denoted P and Q , with respective probabilities p and $q = 1 - p$, then the probability of obtaining n occurrences of outcome P in N observations is

$$P_N(n) = \frac{N!}{n!(N-n)!} p^n q^{N-n}. \quad (5.27)$$

Thus, making use of Equation (5.21), the mean number of occurrences of outcome P in N observations is given by

$$\langle n \rangle = \sum_{n=0, N} P_N(n) n = \sum_{n=0, N} \frac{N!}{n!(N-n)!} p^n q^{N-n} n. \quad (5.28)$$

We can see that if the final factor n were absent on the right-hand side of the previous expression then it would just reduce to the binomial expansion, which we know how to sum. [See Equation (5.16).] We can take advantage of this fact using a rather elegant mathematical sleight of hand. Observe that because

$$n p^n \equiv p \frac{\partial}{\partial p} p^n, \quad (5.29)$$

the previous summation can be rewritten as

$$\sum_{n=0, N} \frac{N!}{n!(N-n)!} p^n q^{N-n} n \equiv p \frac{\partial}{\partial p} \left[\sum_{n=0, N} \frac{N!}{n!(N-n)!} p^n q^{N-n} \right]. \quad (5.30)$$

The term in square brackets is now the familiar binomial expansion, and can be written more succinctly as $(p + q)^N$. Thus,

$$\sum_{n=0,N} \frac{N!}{n!(N-n)!} p^n q^{N-n} n = p \frac{\partial}{\partial p} (p + q)^N = p N (p + q)^{N-1}. \quad (5.31)$$

However, $p + q = 1$ for the case in hand [see Equation (5.9)], so

$$\langle n \rangle = N p. \quad (5.32)$$

In fact, we could have guessed the previous result. By definition, the probability, p , is the number of occurrences of the outcome P divided by the number of observations, in the limit as the number of observations goes to infinity:

$$p = \lim_{N \rightarrow \infty} \frac{n}{N}. \quad (5.33)$$

[See Equation (5.1).] If we think carefully, however, we can appreciate that taking the limit as the number of observations goes to infinity is equivalent to taking the mean value, so that

$$p = \left\langle \frac{n}{N} \right\rangle = \frac{\langle n \rangle}{N}. \quad (5.34)$$

But, this is just a simple rearrangement of Equation (5.32).

Let us now calculate the variance of n . Recall, from Equation (5.25), that

$$\langle (\Delta n)^2 \rangle = \langle n^2 \rangle - \langle n \rangle^2 \quad (5.35)$$

We already know $\langle n \rangle$, so we just need to calculate $\langle n^2 \rangle$. This average is written

$$\langle n^2 \rangle = \sum_{n=0,N} \frac{N!}{n!(N-n)!} p^n q^{N-n} n^2. \quad (5.36)$$

The sum can be evaluated using a simple extension of the mathematical trick that we used previously to evaluate $\langle n \rangle$. Because

$$n^2 p^n \equiv \left(p \frac{\partial}{\partial p} \right)^2 p^n, \quad (5.37)$$

we can write

$$\begin{aligned} \sum_{n=0,N} \frac{N!}{n!(N-n)!} p^n q^{N-n} n^2 &\equiv \left(p \frac{\partial}{\partial p} \right)^2 \sum_{n=0,N} \frac{N!}{n!(N-n)!} p^n q^{N-n} \\ &= \left(p \frac{\partial}{\partial p} \right)^2 (p + q)^N \\ &= \left(p \frac{\partial}{\partial p} \right) [p N (p + q)^{N-1}] \end{aligned}$$

$$= p [N(p+q)^{N-1} + pN(N-1)(p+q)^{N-2}]. \quad (5.38)$$

Using $p+q=1$, we obtain

$$\begin{aligned} \langle n^2 \rangle &= p [N + pN(N-1)] = Np(1 + pN - p) \\ &= (Np)^2 + Npq = \langle n \rangle^2 + Npq, \end{aligned} \quad (5.39)$$

because $\langle n \rangle = Np$. [See Equation (5.32).] It follows that the variance of n is given by

$$\langle (\Delta n)^2 \rangle = \langle n^2 \rangle - \langle n \rangle^2 = Npq. \quad (5.40)$$

The standard deviation of n is the square root of the variance [see Equation (5.26)], so that

$$\sigma_n = \sqrt{Npq}. \quad (5.41)$$

Now, the standard deviation is essentially the width of the range of probable values over which n is distributed around its mean value, $\langle n \rangle$. The relative width of the distribution is characterized by

$$\frac{\sigma_n}{\langle n \rangle} = \frac{\sqrt{Npq}}{Np} = \sqrt{\frac{q}{p}} \frac{1}{\sqrt{N}}. \quad (5.42)$$

It is clear, from the previous formula, that the relative width decreases with increasing N like $N^{-1/2}$. So, the greater the number of observations, the more likely it is that an observation of n will yield a result that is relatively close to the mean value, $\langle n \rangle$.

5.1.5 Random Walk

The so-called *random walk* is a stochastic process that governs, for example, the path traced by a molecule as it travels through a liquid or a gas, while constantly colliding with the other molecules in the medium. (See Section 5.3.9.)

Consider a random walk in one dimension. Suppose that a molecule takes steps of equal length l along the x -axis. Suppose, further, that the steps are taken to the left (i.e., in the negative x -direction) or to the right, at random, with equal probabilities. Let x_n be the molecule's x coordinate after n steps. It is assumed that $x_0 = 0$. In other words, the molecule is initially at the origin. We can write

$$x_n = x_{n-1} \pm l. \quad (5.43)$$

Hence,

$$x_n^2 = x_{n-1}^2 \pm 2x_n l + l^2, \quad (5.44)$$

which implies that

$$\langle x_n^2 \rangle = \langle x_{n-1}^2 \rangle + l^2. \quad (5.45)$$

Thus, by induction, after N steps, we obtain

$$\langle x^2 \rangle = Nl^2 \quad (5.46)$$

Suppose that the steps are taken at a mean frequency f . It follows that $N = ft$, where $x = 0$ at $t = 0$. Hence,

$$\langle x^2 \rangle = 2Dt, \quad (5.47)$$

where

$$D = \frac{1}{2} fl^2 \quad (5.48)$$

is known as the *diffusivity*. According to Equation (5.47), the molecule's mean square distance from its starting point grows linearly in time. This type of motion is known as *diffusion*. (See Section 5.3.9.)

Consider a random walk in three dimensions. Let \mathbf{r} be the displacement of our molecule from the origin (which is its starting point). Suppose that the molecule takes steps of uniform length l , in a random direction, f times a second. Let \mathbf{l} be the displacement associated with a given step. Let \mathbf{r}_n be the molecule's displacement after n steps. We can write

$$\mathbf{r}_n = \mathbf{r}_{n-1} + \mathbf{l}. \quad (5.49)$$

Thus,

$$r_n^2 = (\mathbf{r}_{n-1} + \mathbf{l}) \cdot (\mathbf{r}_{n-1} + \mathbf{l}) = r_{n-1}^2 + 2\mathbf{r}_{n-1} \cdot \mathbf{l} + l^2. \quad (5.50)$$

However, if \mathbf{l} is in a random direction then $\langle \mathbf{r}_{n-1} \cdot \mathbf{l} \rangle = 0$, because the cosine of the angle subtended between \mathbf{r}_{n-1} and \mathbf{l} is just as likely to be positive as to be negative. Hence, the average of the previous equation yields

$$\langle r_n^2 \rangle = \langle r_{n-1}^2 \rangle + l^2. \quad (5.51)$$

By induction, after N steps, we obtain

$$\langle r^2 \rangle = Nl^2, \quad (5.52)$$

which implies that

$$\langle r^2 \rangle = 2Dt, \quad (5.53)$$

where $D = (1/2) fl^2$. Thus, the motion of the molecule is again diffusive in nature.

5.1.6 Continuous Probability Distribution

Consider some physical system A . Suppose that a measurement of a particular physical property of this system, x , can result in a continuous range of different outcomes such that $-\infty < x < \infty$. Now, we would expect the probability that a measurement of x yields a result in the range x to $x + dx$ to be proportional to dx , in the limit that $dx \rightarrow 0$. (See Section 5.1.7.) Hence, we can define the *probability density*, $P(x)$, such that the probability of a measurement of x yielding a result in the range x to $x + dx$ is $P(x)dx$. A simple extension of the result (5.3) yields the normalization condition,

$$\int_{-\infty}^{\infty} P(x) dx = 1. \quad (5.54)$$

It follows, from a straightforward extension of the results in Section 5.1.3 that the mean value of x is

$$\langle x \rangle = \int_{-\infty}^{\infty} P(x) x dx, \quad (5.55)$$

the mean value of x^2 is

$$\langle x^2 \rangle = \int_{-\infty}^{\infty} P(x) x^2 dx, \quad (5.56)$$

and the variance of x is again

$$\langle (\Delta x)^2 \rangle = \langle x^2 \rangle - \langle x \rangle^2. \quad (5.57)$$

If $X(x)$ is some function of x then

$$\langle X \rangle = \int_{-\infty}^{\infty} P(x) X(x) dx. \quad (5.58)$$

Moreover, if $X(x)$ and $Y(x)$ are independent functions of x then

$$\langle X + Y \rangle = \int_{-\infty}^{\infty} P(x) [X(x) + Y(x)] dx = \int_{-\infty}^{\infty} P(x) X(x) dx + \int_{-\infty}^{\infty} P(x) Y(x) dx = \langle X \rangle + \langle Y \rangle. \quad (5.59)$$

Finally, in some situations it is convenient to use a probability density, $P(x)$, that does not satisfy the normalization condition (5.54). In such situations,

$$\langle X \rangle = \frac{\int_{-\infty}^{\infty} P(x) X(x) dx}{\int_{-\infty}^{\infty} P(x) dx}. \quad (5.60)$$

5.1.7 Gaussian Probability Distribution

Consider a very large number of observations, $N \gg 1$, made on a system with two possible outcomes. (See Sections 5.1.2 and 5.1.4.) Suppose that the probability of outcome P is sufficiently large that the average number of occurrences after N observations is much greater than unity; that is,

$$\langle n \rangle = N p \gg 1. \quad (5.61)$$

In this limit, the standard deviation of n is also much greater than unity,

$$\sigma_n = \sqrt{N p q} \gg 1, \quad (5.62)$$

implying that there are very many probable values of n scattered about the mean value, $\langle n \rangle$. This suggests that the probability of obtaining n occurrences of outcome P does not change significantly in going from one possible value of n to an adjacent value. In other words,

$$\frac{|P_N(n+1) - P_N(n)|}{P_N(n)} \ll 1. \quad (5.63)$$

In this situation, it is useful to regard the probability as a smooth function of n . Let n now be a continuous variable that is interpreted as the number of occurrences of outcome P (after N observations) whenever it takes on a positive integer value. The probability that n lies between n and $n + dn$ is defined

$$P(n, n + dn) = P(n) dn, \quad (5.64)$$

where $P(n)$ is a probability density (see Section 5.1.6), and is independent of dn . The probability can be written in this form because $P(n, n + dn)$ can always be expanded as a Taylor series in dn , and must go to zero as $dn \rightarrow 0$. We can write

$$\int_{n-1/2}^{n+1/2} P(n) dn = P_N(n), \quad (5.65)$$

which is equivalent to smearing out the discrete probability $P_N(n)$ over the range $n \pm 1/2$. Given Equations (5.27) and (5.63), the previous relation can be approximated as

$$P(n) \simeq P_N(n) = \frac{N!}{n! (N - n)!} p^n q^{N-n}. \quad (5.66)$$

For large N , the relative width of the probability distribution function is small; that is,

$$\frac{\sigma_n}{\langle n \rangle} = \sqrt{\frac{q}{p}} \frac{1}{\sqrt{N}} \ll 1. \quad (5.67)$$

This suggests that $P(n)$ is strongly peaked around the mean value, $\langle n \rangle$. Suppose that $\ln P(n)$ attains its maximum value at $n = \tilde{n}$ (where we expect $\tilde{n} \sim \langle n \rangle$). Let us Taylor expand $\ln P(n)$ around $n = \tilde{n}$. Note that we are expanding the slowly-varying function $\ln P(n)$, rather than the rapidly-varying function $P(n)$, because the Taylor expansion of $P(n)$ does not converge sufficiently rapidly in the vicinity of $n = \tilde{n}$ to be useful. We can write

$$\ln P(\tilde{n} + \eta) \simeq \ln P(\tilde{n}) + \eta B_1 + \frac{\eta^2}{2} B_2 + \cdots, \quad (5.68)$$

where

$$B_k = \left. \frac{d^k \ln P}{dn^k} \right|_{n=\tilde{n}}. \quad (5.69)$$

By definition,

$$B_1 = 0, \quad (5.70)$$

$$B_2 < 0, \quad (5.71)$$

if $n = \tilde{n}$ corresponds to the maximum value of $\ln P(n)$.

It follows from Equation (5.66) that

$$\ln P = \ln N! - \ln n! - \ln (N - n)! + n \ln p + (N - n) \ln q. \quad (5.72)$$

If n is a large integer, such that $n \gg 1$, then $\ln n!$ is almost a continuous function of n , because $\ln n!$ changes by only a relatively small amount when n is incremented by unity. Hence,

$$\frac{d \ln n!}{dn} \simeq \frac{\ln(n+1)! - \ln n!}{1} = \ln \left[\frac{(n+1)!}{n!} \right] = \ln(n+1), \quad (5.73)$$

giving

$$\frac{d \ln n!}{dn} \simeq \ln n, \quad (5.74)$$

for $n \gg 1$. The integral of this relation

$$\ln n! \simeq n \ln n - n + O(1), \quad (5.75)$$

valid for $n \gg 1$, is called *Stirling's approximation*, after James Stirling, who first obtained it in 1730.

According to Equations (5.69), (5.72), and (5.74),

$$B_1 = -\ln \tilde{n} + \ln(N - \tilde{n}) + \ln p - \ln q. \quad (5.76)$$

Hence, if $B_1 = 0$ then

$$(N - \tilde{n})p = \tilde{n}q, \quad (5.77)$$

giving

$$\tilde{n} = Np = \langle n \rangle, \quad (5.78)$$

because $p + q = 1$. [See Equations (5.9) and (5.32).] Thus, the maximum of $\ln P(n)$ occurs exactly at the mean value of n .

Further differentiation of Equation (5.76) yields [see Equation (5.69)]

$$B_2 = -\frac{1}{\tilde{n}} - \frac{1}{N - \tilde{n}} = -\frac{1}{Np} - \frac{1}{N(1-p)} = -\frac{1}{Npq}, \quad (5.79)$$

because $p + q = 1$. Note that $B_2 < 0$, as required. According to Equation (5.62), the previous relation can also be written

$$B_2 = -\frac{1}{\sigma_n^2}. \quad (5.80)$$

It follows, from the previous analysis, that the Taylor expansion of $\ln P(n)$ can be written

$$\ln P(\langle n \rangle + \eta) \simeq \ln P(\langle n \rangle) - \frac{\eta^2}{2\sigma_n^2} + \dots \quad (5.81)$$

Taking the exponential of both sides, we obtain

$$P(n) \simeq P(\langle n \rangle) \exp \left[-\frac{(n - \langle n \rangle)^2}{2\sigma_n^2} \right]. \quad (5.82)$$

The constant $P(\langle n \rangle)$ is most conveniently fixed by making use of the normalization condition,

$$\int_0^N P(n) dn \simeq 1, \quad (5.83)$$

for a continuous distribution function. [See Equation (5.54). Note that n cannot take a negative value.] Because we only expect $P(n)$ to be significant when n lies in the relatively narrow range $\langle n \rangle \pm \sigma_n^2$, the limits of integration in the previous expression can be replaced by $\pm\infty$ with negligible error. Thus,

$$P(\langle n \rangle) \int_{-\infty}^{\infty} \exp\left[-\frac{(n - \langle n \rangle)^2}{2\sigma_n^2}\right] dn = P(\langle n \rangle) \sqrt{2} \sigma_n \int_{-\infty}^{\infty} e^{-x^2} dx \simeq 1. \quad (5.84)$$

As is well known,

$$\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}. \quad (5.85)$$

It follows from the normalization condition (5.84) that

$$P(\langle n \rangle) \simeq \frac{1}{\sqrt{2\pi} \sigma_n}. \quad (5.86)$$

Finally, we obtain

$$P(n) \simeq \frac{1}{\sqrt{2\pi} \sigma_n} \exp\left[-\frac{(n - \langle n \rangle)^2}{2\sigma_n^2}\right]. \quad (5.87)$$

This is probability distribution is known as *Gaussian probability distribution*, after the Carl F. Gauss, who discovered in 1809 it while investigating the distribution of errors in measurements. The Gaussian distribution is only valid in the limits $N \gg 1$ and $\langle n \rangle \gg 1$. According to this distribution, at one standard deviation away from the mean value—that is $n = \langle n \rangle \pm \sigma_n$ —the probability density is about 61% of its peak value. At two standard deviations away from the mean value, the probability density is about 13.5% of its peak value. Finally, at three standard deviations away from the mean value, the probability density is only about 1% of its peak value. We conclude that there is very little chance that n lies more than about three standard deviations away from its mean value. In other words, n is almost certain to lie in the relatively narrow range $\langle n \rangle \pm 3\sigma_n$.

Consider the drunken walk discussed at the end of Section 5.1.2. Suppose that the drunken man is equally likely to take a step to the right as to take a step to the left. In other words, $p = q = 1/2$. Thus, according to Equations (5.32) and (5.41),

$$\langle n \rangle = \frac{N}{2}, \quad (5.88)$$

$$\sigma_n = \frac{\sqrt{N}}{2}. \quad (5.89)$$

Equations (5.18) and (5.19) state that the probability of the drunken man taking m net steps to the right after N total steps is

$$P'_N(m) = P_N(n), \quad (5.90)$$

where

$$n = \frac{N + m}{2}. \quad (5.91)$$

In the limit of very many steps, we can treat m and n as continuous variables. Let $P'(m) dm$ be the probability that m lies between m and $m + dm$. Likewise, let $P(n) dn$ be the probability that n lies between n and $n + dn$. It follows that

$$P'(m) dm = P(n) dn, \quad (5.92)$$

where m and n satisfy Equation (5.91). Hence,

$$P'(m) = \frac{1}{2} P\left(\frac{N + m}{2}\right) = \frac{1}{\sqrt{2\pi N}} \exp\left(-\frac{m^2}{2N}\right), \quad (5.93)$$

where use has been made of Equations (5.87), (5.88), (5.89), and (5.91). Suppose that each step is of length l , and that the man takes f steps per second. It follows that the man's displacement from his starting point is $x = ml$. Moreover, $N = ft$. Let $P(x, t) dx$ be the probability that the man's displacement from his starting point after t seconds lies between x and $x + dx$. We have $P(x, t) dx = P'(m) dm$, which implies that $P(x) = P'(m)/l$. Hence, we obtain

$$P(x) = \frac{1}{\sqrt{4\pi Dt}} \exp\left(-\frac{x^2}{4Dt}\right), \quad (5.94)$$

where

$$D = \frac{1}{2} fl^2 \quad (5.95)$$

is the diffusivity. It is easily demonstrated that

$$\langle x^2 \rangle = 2Dt. \quad (5.96)$$

Thus, it is evident from the analysis of Section 5.1.5 that the probability density distribution (5.94) corresponds to that of a random walk in one dimension. Equation (5.94) can also be thought of as describing the diffusion of probability density along the x -axis. (See Section 5.3.9.)

5.2 Ideal Gas

5.2.1 Ideal Gas Law

An ideal gas consists of molecules of negligible spatial extent that do not exert forces on one another, except when they collide. At sufficiently large temperatures, and sufficiently low mass densities, most gases in nature can be approximated as ideal gases.

According to *Boyle's law*, which is an experimental result that was first reported by Robert Boyle in 1660, the pressure of an ideal gas is inversely proportional to its volume, at fixed temperature. According to *Charles's law*, which is another experimental result that was obtained by Jacques Charles in 1787, the volume of an ideal gas is proportional to its absolute temperature, at

fixed pressure. Finally, according to *Avogadro's law*, which was first proposed by Amedeo Avogadro in 1812, equal volumes of all ideal gases, at the same temperature and pressure, contain the same number of molecules. These three laws imply that an ideal gas is governed by the following equation of state:

$$pV = \nu RT. \quad (5.97)$$

Here, p is the gas pressure, V the volume, T the absolute temperature, ν the number of moles of molecules in the gas, and

$$R = 8.3145 \text{ J K}^{-1} \text{ mol}^{-1} \quad (5.98)$$

is a constant of proportionality known as the *ideal gas constant*. Equation (5.97) is called the *ideal gas law*. Note that V and ν are *extensive* quantities. That is, if we double the size of the system (by combining two identical systems) then we double the values of these quantities. On the other hand, p and T are *intensive* quantities. That is, if we double the size of the system then the values of these quantities are left unchanged.

Absolute temperature is measured in degrees kelvin (K) on a scale in which absolute zero (i.e., the lowest possible temperature) is 0 K, and the triple point of water (i.e., the unique temperature at which all three phases of water coexist) is 273.16 K.

One mole of molecules contains *Avogadro's number* of molecules; that is,

$$N_A = 6.0221 \times 10^{23} \quad (5.99)$$

molecules. Finally, the *Boltzmann constant*, k_B , is defined

$$k_B = \frac{R}{N_A} = 1.3806 \times 10^{-23} \text{ J K}^{-1}. \quad (5.100)$$

5.2.2 First Law of Thermodynamics

Let U be the *internal energy* of an ideal gas. Internal energy is the energy that the gas possesses by virtue of the random motions of its constituent molecules. Consider a process by which an infinitesimal amount of heat, dQ , is absorbed by the gas, and an infinitesimal amount of work, dW , is performed on the gas. According to the *first law of thermodynamics*, which is a statement of energy conservation that was first explicitly formulated by Rudolf Clausius in 1850,

$$dU = dQ + dW. \quad (5.101)$$

In reality, dQ cannot be directly measured, and is, instead, inferred to be the difference between the change in the gas's internal energy and the work performed on the gas, both of which can be directly measured, according to the previous equation.

Consider an ideal gas in a cylindrical container of cross-sectional area A . Suppose that the top of the container is a movable piston, and that the gas pushes the piston upward a distance dx . Now, from the definition of pressure, the gas exerts a force pA on the piston. Thus, the gas does work $pA dx$ on the piston. (See Section 1.3.2.) Hence, the work done on the gas is $dW = -pA dx = -p dV$, where $dV = A dx$ is the change in volume of the gas. This is a general result. Hence, Equation (5.101) becomes

$$dU = dQ - p dV. \quad (5.102)$$

5.2.3 Specific Heat Capacity

Suppose that we add an amount of heat dQ to an ideal gas causing its temperature to rise by dT . The *specific heat capacity* of the gas is defined

$$C = \frac{dQ}{dT}. \quad (5.103)$$

In fact, an ideal gas possesses a number of different specific heat capacities depending on what is held constant as heat is added to the system. Suppose that the volume is held constant. It follows from Equation (5.102) that $dQ = dU$. Hence, the *specific heat capacity at constant volume* of the gas is

$$C_V = \left(\frac{\partial Q}{\partial T} \right)_V = \left(\frac{\partial U}{\partial T} \right)_V. \quad (5.104)$$

However, according to *Joule's second law*, which was established experimentally by James Joule in 1843, the internal energy of an ideal gas depends on its temperature alone, and is independent of the volume or pressure. We also expect C_V to be an extensive quantity. It follows that

$$dU = \nu c_V(T) dT, \quad (5.105)$$

where $c_V = C_V/\nu$ is termed the *molar specific heat capacity at constant volume*, and is an intensive quantity. In fact, c_V is constant for an idea gas. For a monatomic gas,

$$c_V = \frac{3}{2} R, \quad (5.106)$$

whereas for a diatomic gas,

$$c_V = \frac{5}{2} R. \quad (5.107)$$

(See Sections 5.3.6 and 5.5.8.) In both cases, Equation (5.105) can be integrated to give

$$U(T) = \nu c_V T. \quad (5.108)$$

Consider the specific heat capacity of an ideal gas at constant pressure. Making use of Equations (5.102) and (5.105),

$$dQ = \nu c_V dT + p dV. \quad (5.109)$$

However, at constant pressure, the ideal gas law, (5.97), yields

$$p dV = \nu R dT. \quad (5.110)$$

The previous two equations give

$$dQ = \nu (c_V + R) dT. \quad (5.111)$$

Now, the *molar specific heat capacity at constant pressure* of an ideal gas is

$$c_p = \frac{1}{\nu} \left(\frac{\partial Q}{\partial T} \right)_p. \quad (5.112)$$

Hence, we deduce that

$$c_p = c_v + R. \quad (5.113)$$

Note that the specific heat capacity at constant pressure is greater than that at constant volume, because some of the heat energy added to a gas held at constant pressure is consumed by the work that the gas does on its surroundings, in order to expand its volume slightly, and, therefore, does not lead to an increase in the internal energy (i.e., temperature) of the gas. On the other hand, for a gas held at constant volume, all of the added heat energy goes to increase its internal energy.

5.2.4 Isothermal and Adiabatic Expansion

Suppose that the temperature of an ideal gas is held constant by keeping the gas in thermal contact with a heat reservoir. If the gas is allowed to expand quasi-statically under these so-called *isothermal* conditions then the ideal gas law, (5.97), tells us that

$$pV = \text{constant}. \quad (5.114)$$

This result is known as the *isothermal gas law*.

Suppose, now, that the gas is thermally isolated from its surroundings. If the gas is allowed to expand quasi-statically under these so-called *adiabatic* conditions then it does work on its environment, and, hence, its internal energy is reduced, and its temperature decreases. Let us calculate the relationship between the pressure and volume of the gas during adiabatic expansion. According to Equation (5.109),

$$dQ = \nu c_v dT + p dV = 0, \quad (5.115)$$

in an adiabatic process (in which no heat is absorbed). The ideal gas law, (5.97), can be differentiated, yielding

$$p dV + V dp = \nu R dT. \quad (5.116)$$

The temperature increment, dT , can be eliminated between the previous two expressions to give

$$0 = \frac{c_v}{R} (p dV + V dp) + p dV = \left(\frac{c_v}{R} + 1 \right) p dV + \frac{c_v}{R} V dp, \quad (5.117)$$

which reduces to

$$(c_v + R) p dV + c_v V dp = 0. \quad (5.118)$$

Dividing through by $c_v p V$ yields

$$\gamma \frac{dV}{V} + \frac{dp}{p} = 0, \quad (5.119)$$

where

$$\gamma \equiv \frac{c_p}{c_v} = \frac{c_v + R}{c_v} \quad (5.120)$$

is termed the *ratio of specific heats*. [See Equation (5.113).] Given that c_v is a constant in an ideal gas, the ratio of specific heats, γ , is also a constant. In fact, Equations (5.106), (5.107), and the previous equation, imply that

$$\gamma = \frac{5}{3} \quad (5.121)$$

for a monatomic gas, and

$$\gamma = \frac{7}{5} \quad (5.122)$$

for a diatomic gas.

Because γ is a constant for an ideal gas, we can integrate Equation (5.119) to give

$$\gamma \ln V + \ln p = \text{constant}, \quad (5.123)$$

or

$$p V^\gamma = \text{constant}. \quad (5.124)$$

This result is known as the *adiabatic gas law*. It is straightforward to obtain analogous relationships between V and T , and between p and T , during adiabatic expansion or contraction. In fact, because $p = \nu R T / V$, the previous formula also implies that

$$T V^{\gamma-1} = \text{constant}, \quad (5.125)$$

and

$$p^{1-\gamma} T^\gamma = \text{constant}. \quad (5.126)$$

Equations (5.124)–(5.126) are all completely equivalent.

5.2.5 Hydrostatic Equilibrium of Atmosphere

The gas that we are most familiar with in everyday life is, of course, the Earth's atmosphere. It turns out that we can use the isothermal and adiabatic gas laws to explain most of the observed features of the atmosphere.

Let us, first of all, consider the hydrostatic equilibrium of the atmosphere. Consider a thin vertical slice of the atmosphere, of cross-sectional area A , that starts at height z above ground level, and extends to height $z + dz$. The upward force exerted on this slice by the gas below it is $p(z)A$, where $p(z)$ is the pressure at height z . Likewise, the downward force exerted by the gas above the slice is $p(z + dz)A$. The net upward force is $[p(z) - p(z + dz)]A$. In equilibrium, this upward force must be balanced by the downward force due to the weight of the slice, which is $\rho A dz g$, where ρ is the mass density of the gas, and g the acceleration due to gravity. It follows that the force balance condition can be written

$$[p(z) - p(z + dz)]A = \rho A dz g, \quad (5.127)$$

which reduces to

$$\frac{dp}{dz} = -\rho g. \quad (5.128)$$

This result is known as the *equation of hydrostatic equilibrium* for the atmosphere.

We can express the mass density of a gas in the following form,

$$\rho = \frac{\nu \mu}{V}, \quad (5.129)$$

where μ is the *molecular weight* of the gas, and is equal to the mass of one mole of gas particles. For instance, the molecular weight of nitrogen gas is 28×10^{-3} kg. The previous formula for the mass density of a gas, combined with the ideal gas law, $pV = \nu RT$, yields

$$\rho = \frac{p\mu}{RT}. \quad (5.130)$$

It follows that the equation of hydrostatic equilibrium can be rewritten

$$\frac{dp}{p} = -\frac{\mu g}{RT} dz. \quad (5.131)$$

5.2.6 Isothermal Atmosphere

As a first approximation, let us assume that the temperature of the atmosphere is uniform. In such an *isothermal atmosphere*, we can directly integrate the previous equation to give

$$p = p_0 \exp\left(-\frac{z}{z_0}\right). \quad (5.132)$$

Here, p_0 is the pressure at ground level ($z = 0$), which is generally about 1 bar (10^5 Nm⁻² in SI units). The quantity

$$z_0 = \frac{RT}{\mu g} \quad (5.133)$$

is known as the *isothermal scale-height* of the atmosphere. At ground level, the atmospheric temperature is, on average, about 15°C, which is 288K on the absolute scale. The mean molecular weight of air at sea level is 29×10^{-3} kg (i.e., the molecular weight of a gas made up of 78% nitrogen, 21% oxygen, and 1% argon). The mean acceleration due to gravity is 9.81 m s⁻² at ground level. Also, the molar ideal gas constant is 8.314 joules/mole/degree. Combining all of this information, the isothermal scale-height of the atmosphere comes out to be about 8.4 kilometers.

We have discovered that, in an isothermal atmosphere, the pressure decreases exponentially with increasing height. Because the temperature is assumed to be constant, and $\rho \propto p/T$ [see Equation (5.130)], it follows that the density also decreases exponentially with the same scale-height as the pressure. According to Equation (5.132), the pressure, or the density, of the atmosphere decreases by a factor 10 every $\ln 10 z_0$, or 19.3 kilometers, increase in altitude above sea level. Clearly, the effective height of the atmosphere is very small compared to the Earth's radius, which is about 6,400 kilometers. In other words, the atmosphere constitutes a relatively thin layer covering the surface of the Earth. Incidentally, this justifies our neglect of the decrease of g with increasing altitude.

One of the highest points in the United States of America is the peak of Mount Elbert in Colorado. This peak lies 14,432 feet, or about 4.4 kilometers, above sea level. At this altitude, Equation (5.132) predicts that the air pressure should be about 0.6 atmospheres. Surprisingly enough, after a few days acclimatization, humans can survive quite comfortably at this sort of pressure. In the highest inhabited regions of the Andes and Tibet, the air pressure falls to about 0.5 atmospheres. Humans can just about survive at such pressures. However, humans cannot survive

for any extended period in air pressures below half an atmosphere. This sets an upper limit on the altitude of permanent human habitation, which is about 19,000 feet, or 5.8 kilometers, above sea level.

The highest point in the world is, of course, the peak of Mount Everest in Nepal. This peak lies at an altitude of 29,028 feet, or 8.85 kilometers, above sea level, where we expect the air pressure to be a mere 0.35 atmospheres. This explains why Mount Everest was only conquered after lightweight portable oxygen cylinders were invented. Admittedly, some climbers have subsequently ascended Mount Everest without the aid of additional oxygen, but this is a very foolhardy venture, because, above 19,000 feet, the climbers are slowly dying.

Commercial airliners fly at a cruising altitude of 32,000 feet. At this altitude, we expect the air pressure to be only 0.3 atmospheres, which explains why airline cabins are pressurized. In fact, the cabins are only pressurized to 0.85 atmospheres (which accounts for the “popping” of passengers ears during air travel). The reason for this partial pressurization is quite simple. At 32,000 feet, the pressure difference between the air in the cabin and the air outside the aircraft is about half an atmosphere. Clearly, the walls of the cabin must be strong enough to support this pressure difference, which implies that they must be of a certain thickness, and, hence, that the aircraft must be of a certain weight. If the cabin were fully pressurized then the pressure difference at cruising altitude would increase by about 30%, which implies that the cabin walls would have to be much thicker, and, hence, the aircraft would have to be substantially heavier. So, a fully pressurized aircraft would be more comfortable to fly in (because your ears would not “pop”), but it would also be far less economical to operate.

5.2.7 Adiabatic Atmosphere

Of course, we know that the atmosphere is not isothermal. In fact, air temperature falls quite noticeably with increasing altitude. In ski resorts, the general rule of thumb is that the temperature drops by about 1 degree per 100 meters increase in altitude. Many people cannot understand why the atmosphere gets colder with increasing height. They reason that because higher altitudes are closer to the Sun they ought to be hotter. In fact, the explanation is quite simple. It depends on three important properties of air. The first property is that air is transparent to most, but by no means all, of the electromagnetic spectrum. In particular, most infrared radiation, which carries heat energy, passes straight through the lower atmosphere, and heats the ground. In other words, the lower atmosphere is heated from below, not from above. The second important property of air is that it is constantly in motion. In fact, the lower 20 kilometers of the atmosphere (the so-called *troposphere*) are fairly thoroughly mixed. You might think that this would imply that the atmosphere is isothermal. However, this is not the case because of the final important property of air; namely, it is a very poor conductor of heat. (See Section 5.3.10.) This, of course, is why woolly sweaters work; they trap a layer of air close to the body, and, because air is such a poor conductor of heat, you stay warm.

Imagine a packet of air that is swirling around in the atmosphere. We would expect it to always remain at the same pressure as its surroundings, otherwise it would be mechanically unstable. It is also plausible that the packet moves around too quickly to effectively exchange heat with its surroundings, because air is very a poor heat conductor, and heat flow is consequently quite a slow

process. So, to a first approximation, the air in the packet is adiabatic. In a steady-state atmosphere, we expect that, as the packet moves upwards, expands due to the reduced pressure, and cools adiabatically, its temperature always remains the same as that of its immediate surroundings. This implies that we can use the adiabatic gas law to characterize the cooling of the atmosphere with increasing altitude. In this particular case, the most useful manifestation of the adiabatic law is

$$p^{1-\gamma} T^\gamma = \text{constant}, \quad (5.134)$$

giving

$$\frac{dp}{p} = \frac{\gamma}{\gamma-1} \frac{dT}{T}. \quad (5.135)$$

Combining the previous expression with the equation of hydrostatic equilibrium, (5.131), we obtain

$$\frac{\gamma}{\gamma-1} \frac{dT}{T} = -\frac{\mu g}{RT} dz, \quad (5.136)$$

or

$$\frac{dT}{dz} = -\frac{\gamma-1}{\gamma} \frac{\mu g}{R}. \quad (5.137)$$

Now, the ratio of specific heats for air (which is effectively a diatomic gas) is about 1.4. [See Equation (5.122).] Hence, given that $\mu = 29 \times 10^{-3} \text{ kg}$ and $g = 9.81 \text{ m s}^{-2}$, we deduce, from the previous expression, that the temperature of the atmosphere decreases with increasing height at a constant rate of 9.8°C per kilometer. This value is called the (dry) *adiabatic lapse rate* of the atmosphere. Our calculation accords well with the “1 degree colder per 100 meters higher” rule of thumb used in ski resorts. The basic reason that air is colder at higher altitudes is that it expands as its pressure decreases with height. It, therefore, does work on its environment, without absorbing any heat (because of its low thermal conductivity), so its internal energy, and, hence, its temperature decreases.

According to the adiabatic lapse rate calculated previously, the air temperature at the cruising altitude of airliners (32,000 feet) should be about -80°C (assuming a sea level temperature of 15°C). In fact, this is somewhat of an underestimate. A more realistic value is about -60°C . The explanation for this discrepancy is the presence of water vapor in the atmosphere. As air rises, expands, and cools, water vapor condenses out, releasing latent heat, which prevents the temperature from falling as rapidly with height as the adiabatic lapse rate would predict. In fact, in the tropics, where the air humidity is very high, the lapse rate of the atmosphere (i.e., the rate of decrease of temperature with altitude) is significantly less than the adiabatic value. The adiabatic lapse rate is only observed when the humidity is low. This is the case in deserts, in the arctic (where water vapor is frozen out of the atmosphere), and, of course, in ski resorts.

Suppose that the lapse rate of the atmosphere differs from the adiabatic value. Let us ignore the complication of water vapor, and assume that the atmosphere is dry. Consider a packet of air that moves slightly upwards from its equilibrium height. The temperature of the packet will decrease with altitude according to the adiabatic lapse rate, because its expansion is adiabatic. We shall assume that the packet always maintains pressure balance with its surroundings. It follows that because $\rho T \propto p$, according to the ideal gas law,

$$(\rho T)_{\text{packet}} = (\rho T)_{\text{atmosphere}}. \quad (5.138)$$

If the atmospheric lapse rate is less than the adiabatic value then $T_{\text{atmosphere}} > T_{\text{packet}}$ implying that $\rho_{\text{packet}} > \rho_{\text{atmosphere}}$. So, the packet will be denser than its immediate surroundings, and will, therefore, tend to fall back to its original height. Clearly, an atmosphere whose lapse rate is less than the adiabatic value is vertically stable. On the other hand, if the atmospheric lapse rate exceeds the adiabatic value then, after rising a little way, the packet will be less dense than its immediate surroundings, and will, therefore, continue to rise due to buoyancy effects. Clearly, an atmosphere whose lapse rate is greater than the adiabatic value is vertically unstable. This effect is of great importance in meteorology. The normal stable state of the atmosphere is for the lapse rate to be slightly less than the adiabatic value. Occasionally, however, the lapse rate exceeds the adiabatic value, and this is always associated with extremely disturbed weather patterns.

Let us consider the temperature, pressure, and density profiles in an adiabatic atmosphere. We can directly integrate Equation (5.137) to give

$$T = T_0 \left(1 - \frac{\gamma - 1}{\gamma} \frac{z}{z_0} \right), \quad (5.139)$$

where T_0 is the ground-level temperature, and

$$z_0 = \frac{RT_0}{\mu g} \quad (5.140)$$

the isothermal scale-height calculated using this temperature. The pressure profile is easily calculated from the adiabatic gas law $p^{1-\gamma} T^\gamma = \text{constant}$, or $p \propto T^{\gamma/(\gamma-1)}$. It follows that

$$p = p_0 \left(1 - \frac{\gamma - 1}{\gamma} \frac{z}{z_0} \right)^{\gamma/(\gamma-1)}. \quad (5.141)$$

Consider the limit $\gamma \rightarrow 1$. In this limit, Equation (5.139) yields T independent of height (i.e., the atmosphere becomes isothermal). We can evaluate Equation (5.141) in the limit as $\gamma \rightarrow 1$ using the mathematical identity

$$\lim_{m \rightarrow 0} (1 + mx)^{1/m} \equiv \exp(x). \quad (5.142)$$

We obtain

$$p = p_0 \exp\left(-\frac{z}{z_0}\right), \quad (5.143)$$

which, not surprisingly, is the predicted pressure variation in an isothermal atmosphere. In reality, the ratio of specific heats of the atmosphere is not unity, but is about 1.4 (i.e., the ratio for diatomic gases), which implies that in the real atmosphere

$$p = p_0 \left(1 - \frac{z}{3.5 z_0} \right)^{3.5}. \quad (5.144)$$

In fact, this formula gives very similar results to the isothermal formula, Equation (5.143), for heights below one scale-height (i.e., $z < z_0$). For heights above one scale-height, the isothermal formula tends to predict too high a pressure. See Figure 5.1. So, in an adiabatic atmosphere, the

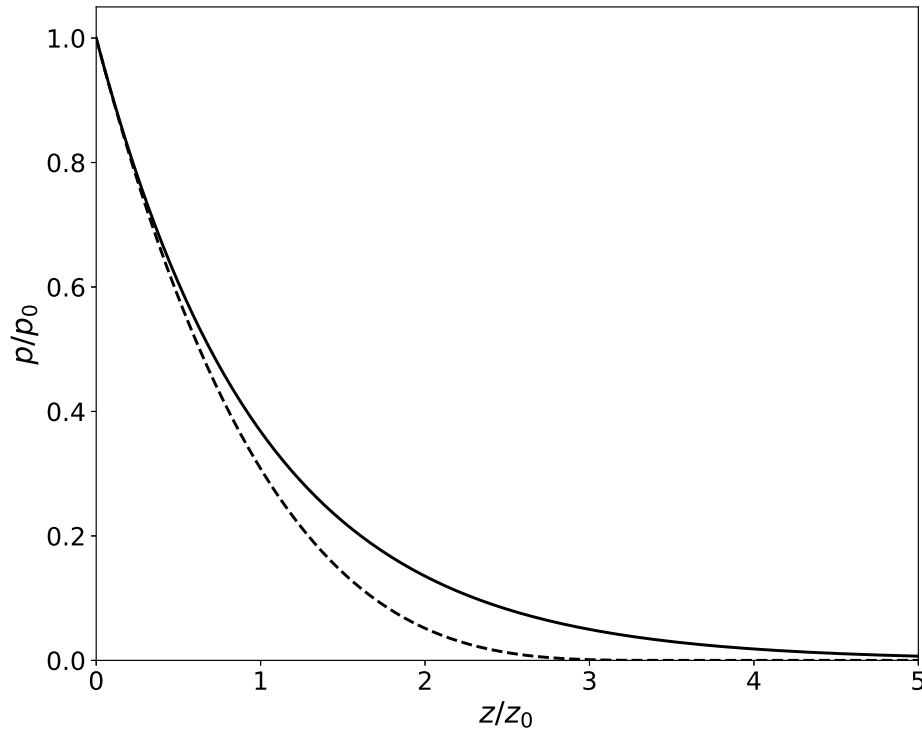


Figure 5.1: The solid curve shows the variation of pressure (normalized to the pressure at ground level) with altitude (normalized to the isothermal scale-height at ground level) in an isothermal atmosphere. The dashed curve shows the variation of pressure with altitude in an adiabatic atmosphere.

pressure falls off more quickly with altitude than in an isothermal atmosphere, but this effect is only noticeable at pressures significantly below one atmosphere. In fact, the isothermal formula is a fairly good approximation below altitudes of about 10 kilometers. Because $\rho \propto p/T$, the variation of density with height is

$$\rho = \rho_0 \left(1 - \frac{\gamma - 1}{\gamma} \frac{z}{z_0}\right)^{1/(\gamma-1)} = \rho_0 \left(1 - \frac{z}{3.5 z_0}\right)^{2.5}, \quad (5.145)$$

where ρ_0 is the density at ground level. Thus, the density falls off more rapidly with altitude than the temperature, but less rapidly than the pressure.

Note that an adiabatic atmosphere has a sharp upper boundary. Above height $z_1 = [\gamma/(\gamma-1)] z_0$, the temperature, pressure, and density are all zero. In other words, there is no atmosphere. For real air, with $\gamma = 1.4$, the upper boundary of an adiabatic atmosphere lies at height $z_1 \approx 3.5 z_0 \approx 29.4$ kilometers above sea level. This behavior is quite different to that of an isothermal atmosphere, which has a diffuse upper boundary. In reality, there is no sharp upper boundary to the atmosphere.

The adiabatic gas law does not apply above about 20 kilometers (i.e., in the *stratosphere*) because, at these altitudes, the air is no longer strongly mixed. Thus, in the stratosphere, the pressure falls off exponentially with increasing height.

5.2.8 Bulk Modulus

The *bulk modulus* of an ideal gas is a measure of its resistance to bulk compression, and is defined

$$\kappa = -V \left(\frac{\partial p}{\partial V} \right). \quad (5.146)$$

In fact, an ideal gas possesses a number of different bulk moduli depending on what is held constant as the pressure is varied. The two most important bulk moduli are the *isothermal bulk modulus*,

$$\kappa_T = -V \left(\frac{\partial p}{\partial V} \right)_T, \quad (5.147)$$

and the *isentropic bulk modulus*,

$$\kappa_S = -V \left(\frac{\partial p}{\partial V} \right)_S. \quad (5.148)$$

The former describes situations in which the gas undergoes isothermal compression, whereas the latter describes situations in which the gas undergoes adiabatic compression. (Note that S actually denotes *entropy*. However, a gas that undergoes compression at constant entropy is such that no heat is added to the gas during the compression.)

According to the isothermal gas law, (5.114),

$$\ln p + \ln V = \text{constant}, \quad (5.149)$$

so

$$-\frac{\partial \ln p}{\partial \ln V} = -\frac{V}{p} \left(\frac{\partial p}{\partial V} \right)_T = 1, \quad (5.150)$$

which implies that

$$\kappa_T = p. \quad (5.151)$$

According to the adiabatic gas law, (5.124),

$$\ln p + \gamma \ln V = \text{constant}, \quad (5.152)$$

$$-\frac{\partial \ln p}{\partial \ln V} = -\frac{V}{p} \left(\frac{\partial p}{\partial V} \right)_S = \gamma, \quad (5.153)$$

which implies that

$$\kappa_S = \gamma p. \quad (5.154)$$

Note that the isentropic bulk modulus of an ideal gas is greater than its isothermal bulk modulus (because $\gamma > 1$). In other words, an ideal gas resists adiabatic compression to a greater degree than it resists isothermal compression. This is the case because during adiabatic compression the work done on the gas causes its temperature to rise, leading to a greater increase in the pressure than would be obtained if the temperature were held constant.

5.2.9 Sound Waves

A sound wave is a type of longitudinal wave that causes a disturbance in the pressure and density of an ideal gas through which it passes. Consider a plane sound wave propagating in the x -direction. Let $\xi(x, t)$ be the longitudinal displacement of the gas associated with the wave. Consider a slab of gas of cross-sectional area A lying between $x - dx/2$ and $x + dx/2$. The mass of the slab is $\rho A dx$, where ρ is gas's mass density. The slab's equation of longitudinal motion is

$$\rho A dx \frac{\partial^2 \xi}{\partial t^2} = A [-p(x + dx/2) + p(x - dx/2)] = -A \frac{\partial p}{\partial x} dx, \quad (5.155)$$

which gives

$$\rho \frac{\partial^2 \xi}{\partial t^2} = -\frac{\partial p}{\partial x}. \quad (5.156)$$

The change in volume of the slab of gas is

$$\delta V = A [\xi(x + dx/2) - \xi(x - dx/2)] = A \frac{\partial \xi}{\partial x} dx, \quad (5.157)$$

which yields

$$\frac{\delta V}{V} = \frac{\partial \xi}{\partial x}, \quad (5.158)$$

because $V = A dx$. However,

$$\frac{\delta V}{V} = -\frac{\delta p}{\kappa}, \quad (5.159)$$

where κ is the bulk modulus. [See Equation (5.146).] Hence,

$$\frac{\partial \xi}{\partial x} = -\frac{\delta p}{\kappa}. \quad (5.160)$$

Equation (5.156) gives

$$\rho \frac{\partial^2}{\partial t^2} \left(\frac{\partial \xi}{\partial x} \right) = -\frac{\partial^2 \delta p}{\partial x^2}, \quad (5.161)$$

writing $p = p_0 + \delta p(x, t)$, where p_0 is a constant background pressure. The previous two equations can be combined to yield

$$\frac{\partial^2 \delta p}{\partial t^2} = v_s^2 \frac{\partial^2 \delta p}{\partial x^2}, \quad (5.162)$$

where

$$v_s = \left(\frac{\kappa}{\rho} \right)^{1/2}. \quad (5.163)$$

Equation (5.162) is a one-dimensional *wave equation* that has the standard solution

$$\delta p(x, t) = \delta p_0 \cos[k(x - v_s t)], \quad (5.164)$$

where δp_0 and k are constants. The previous solution corresponds to a wave-like disturbance in the gas pressure of amplitude δp_0 , wavenumber $\mathbf{k} = k \mathbf{e}_x$, and phase velocity v_s . In other words, Equation (5.163) specifies the speed of sound in an ideal gas.

It remains to determine whether the compression of the gas associated with the passage of a sound wave is isothermal or isentropic. In fact, because ideal gases are relatively poor conductors of heat (see Section 5.3.10), the period of vibration of a sound wave is generally much shorter than the relaxation time necessary for a small element of the gas to exchange energy with the remainder of the gas by means of heat flow. Hence, the compression of the gas associated with the passage of a sound wave is isentropic. It follows from Equations (5.154) and (5.163) that the speed of sound in an ideal gas is

$$v_s = \left(\frac{\kappa_S}{\rho} \right)^{1/2} = \left(\frac{\gamma P}{\rho} \right)^{1/2}. \quad (5.165)$$

Making use of Equations (5.97) and (5.129), the previous equation becomes

$$v_s = \left(\frac{\gamma R T}{\mu} \right)^{1/2}, \quad (5.166)$$

where μ is the molecular mass. Note that the speed of sound in an ideal gas only depends on the gas temperature, and is independent of the pressure.

It is a good approximation to treat the Earth's atmosphere as an ideal gas. The atmosphere is mostly diatomic, which implies that $\gamma = 1.4$. [See Equation (5.122).] Moreover, the molecular weight of the atmosphere is $\mu = 29 \times 10^{-3}$ kg. (See Section 5.2.6.) Hence, the speed of sound in air at 15° C is 340 m s^{-1} .

5.3 Kinetic Theory

5.3.1 Fundamental Assumptions

The purpose of kinetic theory is to deduce the macroscopic properties of an ideal gas from the motions of its constituent molecules. The fundamental assumptions of kinetic theory are that a gas held in a container consists of a very large number of molecules that are in ceaseless motion. Moreover, these molecules are constantly colliding with one another, and also with the walls of container. Furthermore, the pressure acting on the walls of the container is the resultant of all of the reaction forces as the molecules strike and rebound from the walls.

We can make a number of simplifying assumptions in our exploration of kinetic theory. First, the volume of the molecules is assumed to be negligible. Second, the molecules are assumed not to exert forces on one another, except when they collide. Third, the collisions of the molecules with the walls are assumed to be specular. The first two assumptions merely ensure that we are dealing with an ideal gas.

5.3.2 Molecular Flux

Suppose that the molecules in our gas are equally likely to be moving in any direction, and have a distribution of molecular speeds $F(v)$. (See Section 5.5.9.) In other words, the probability that a given molecule has a speed in the range v to $v+dv$ is $F(v) dv$. Let n be the total number of molecules

per unit volume. Let us calculate how many molecules per unit area, per second, pass through the x - y plane in the direction of increasing z . This quantity, Φ_z , is termed the *molecular flux*.

Let θ and ϕ be standard spherical polar angles. (See Section A.23.) We can write the Cartesian components of the velocity of a given molecule, whose molecular speed is v , as $\mathbf{v} = v(\sin\theta \cos\phi, \sin\theta \sin\phi, \cos\theta)$. If the molecules are equally likely to move in any direction then the number of molecules for which θ lies between θ and $\theta + d\theta$, and ϕ lies between ϕ and $\phi + d\phi$, is proportional to $\sin\theta d\theta d\phi$ (i.e., to the amount of solid angle contained in this range of angles). Thus, the number of molecules for which θ lies between θ and $\theta + d\theta$, and ϕ can take any value in the range 0 to 2π , is proportional to $2\pi \sin\theta d\theta$. Hence, given that there are 4π steradians in a complete solid angle, the fraction of molecules for which θ lies between θ and $\theta + d\theta$ is $g(\theta) d\theta$, where

$$g(\theta) = \frac{1}{2} \sin\theta. \quad (5.167)$$

Consider molecules whose speeds lie between v and $v + dv$. The number of such molecules per unit volume is $n F(v) dv$. The number of such molecules per unit volume whose directions of motion subtend an angle lying between θ and $\theta + d\theta$ with the z -axis is $[n F(v) dv] [g(\theta) d\theta]$. All such molecules for which $-v_z < z < 0$ cross the x - y plane in one second. Thus, the number of such molecules per unit area, per second, that cross the x - y plane is

$$d\Phi_z = [n F(v) dv] [g(\theta) d\theta] [v_z] = [n F(v) dv] \left[\frac{1}{2} \sin\theta d\theta \right] [v \cos\theta]. \quad (5.168)$$

Hence, the net flux of molecules across the x - y plane in the direction of increasing z (i.e., with $0 \leq \theta \leq \pi/2$) is

$$\Phi_z = \frac{1}{2} n \int_0^{\pi/2} \sin\theta \cos\theta d\theta \int_0^\infty F(v) v dv, \quad (5.169)$$

which reduces to

$$\Phi_z = \frac{1}{4} n \langle v \rangle, \quad (5.170)$$

where

$$\langle v \rangle = \int_0^\infty F(v) v dv. \quad (5.171)$$

is the mean molecular speed. (See Sections 5.1.6 and 5.5.9.)

For example, if a low-pressure gas is held in a container, the wall of which contains a small hole of area A , then the number of escaping molecules per second is

$$\dot{N} = \frac{1}{4} n \langle v \rangle A. \quad (5.172)$$

This process of molecular escape is known as *molecular effusion*. (See Section 5.3.13.) It turns out that the previous formula is only accurate if the dimensions of the hole are small compared to the typical distance travelled by a molecule in the gas between collisions (this distance is known as the mean free path; see Section 5.3.8). In the opposite limit, the gas flows through the hole according to the laws of continuum fluid dynamics.

5.3.3 Pressure

Suppose that the x - y plane actually corresponds to a wall of the container. Consider, again, molecules whose speeds lie between v and $v + dv$, and whose directions of motion subtend an angle lying between θ and $\theta + d\theta$ with the z -axis. Each such molecule that encounters the wall bounces off it in a specular fashion, and its z -momentum consequently changes by $2m v_z$, where m is the molecular mass. Thus, the normal reaction force per unit area acting on the wall is

$$dp = [2m v_z] [n F(v) dv] [g(\theta) d\theta] [v_z] = [2m v \cos \theta] [n F(v) dv] \left[\frac{1}{2} \sin \theta d\theta \right] [v \cos \theta]. \quad (5.173)$$

[See Equation (5.167).] Hence, the total pressure exerted on the wall is

$$p = nm \int_0^{\pi/2} \sin \theta \cos^2 \theta d\theta \int_0^\infty F(v) v^2 dv, \quad (5.174)$$

which reduces to

$$p = \frac{1}{3} nm \langle v^2 \rangle, \quad (5.175)$$

where

$$\langle v^2 \rangle = \int_0^\infty F(v) v^2 dv. \quad (5.176)$$

is the mean square molecular speed. (See Section 5.5.9.)

However, we can write

$$n = \frac{\nu N_A}{V}, \quad (5.177)$$

where ν is the number of moles of molecules held inside the container, V is the volume of the container, and N_A is Avogadro's number. Equations (5.175) and (5.177) yield

$$\frac{pV}{\nu} = \frac{2}{3} N_A \langle \mathcal{K}_{\text{trans}} \rangle, \quad (5.178)$$

where

$$\langle \mathcal{K}_{\text{trans}} \rangle = \frac{1}{2} m \langle v^2 \rangle \quad (5.179)$$

is the mean translational kinetic energy of a molecule in the gas. Equation (5.178) is consistent with the ideal gas law, (5.97), provided that

$$\langle \mathcal{K}_{\text{trans}} \rangle = \frac{1}{2} m \langle v^2 \rangle = \frac{3}{2} k_B T, \quad (5.180)$$

where $k_B = R/N_A$ is the Boltzmann constant. [See Equation (5.100).]

5.3.4 Law of Equipartition of Energy

The *law of equipartition of energy* is a result in statistical thermodynamics that states that the mean thermal energy associated with each independent quadratic (i.e., proportional to the square of a coordinate or a momentum component) contribution to the total energy of a system consisting of many particles is $(1/2)k_B T$, where T is the temperature of the system. (See Section 5.5.5.) It turns out, however, that this law only applies if the contribution in question is governed by classical (as opposed to quantum mechanical) physics. (See Section 5.5.6.)

Consider a particular constituent molecule of an ideal gas whose mass is m , and whose velocity is \mathbf{v} . The contribution of molecule's translational kinetic energy to the total energy of the whole gas is

$$\frac{1}{2} m v_x^2 + \frac{1}{2} m v_y^2 + \frac{1}{2} m v_z^2 = \frac{p_x^2}{2m} + \frac{p_y^2}{2m} + \frac{p_z^2}{2m}, \quad (5.181)$$

where $\mathbf{p} = m \mathbf{v}$ is the molecular momentum. It can be seen that the contribution of the molecules's translational kinetic energy to the total energy consists of three terms that are quadratic in a momentum component. Hence, according to the law of equipartition of energy, the mean thermal energy associated with the molecules translational kinetic energy is

$$\langle \mathcal{K}_{\text{trans}} \rangle = \frac{1}{2} k_B T + \frac{1}{2} k_B T + \frac{1}{2} k_B T = \frac{3}{2} k_B T, \quad (5.182)$$

in accordance with Equation (5.180).

5.3.5 Partial Pressure

Suppose that an ideal gas consists of N distinct types of molecule. Let a molecule of type i have a number density n_i , a mass m_i , and a velocity \mathbf{v}_i . If we repeat the analysis of Section 5.3.3, taking into account the different types of molecule, then it is easily shown that the total pressure of the gas is

$$p = \frac{1}{3} \sum_{i=1,N} n_i m_i \langle v_i^2 \rangle. \quad (5.183)$$

However, the law of equipartition of energy (see the previous section) implies that

$$\frac{1}{2} m_i \langle v_i^2 \rangle = \frac{3}{2} k_B T. \quad (5.184)$$

Moreover,

$$n_i = \frac{\nu_i N_A}{V}, \quad (5.185)$$

where ν_i is the number of moles of molecules of type i in the gas, and V is the volume of the gas. [See Equation (5.177).] The previous three equations yield

$$p = \sum_{i=1,N} p_i, \quad (5.186)$$

where

$$p_i V = \nu_i R T. \quad (5.187)$$

We conclude that the total pressure of the gas is the sum of the pressures that a gas of each constituent type of molecule would exert independently. This result is known as *Dalton's law*, after John Dalton, who verified it experimentally in 1802. The quantity p_i is known as the *partial pressure* of type- i molecules. Thus, Dalton's law is equivalent to the statement that the total pressure of an ideal gas is the sum of the partial pressures of the individual gases from which it is composed.

5.3.6 Internal Energy

Consider a monatomic gas such as helium. An individual helium atom can only store energy in its translational motion. As we saw in Section 5.3.4, the mean energy associated with this motion is $(3/2)k_B T$. Hence, the internal energy of a gas consisting of ν moles of helium atoms is

$$U = \nu N_A \frac{3}{2} k_B T = \nu \frac{3}{2} R T. \quad (5.188)$$

However, according to Equation (5.108),

$$U = \nu c_V T, \quad (5.189)$$

where c_V is the molar specific heat capacity of the gas at constant volume. The previous two equations imply that the molar specific heat capacity of a helium gas (or any monatomic gas) is

$$c_V = \frac{3}{2} R, \quad (5.190)$$

in accordance with Equation (5.106).

Consider a diatomic gas such as hydrogen. An individual hydrogen molecule can store energy in its translational motion, but can also store energy in its rotational motion. In principle, a molecule has three principal axes of rotation about which it could rotate. (See Sections 1.7.2 and 1.7.3.) Hence, the net rotational kinetic energy is

$$\frac{L_x^2}{2I_{xx}} + \frac{L_y^2}{2I_{yy}} + \frac{L_z^2}{2I_{zz}}, \quad (5.191)$$

where L_x is the angular momentum about the x -axis, I_{xx} is the principal moment of inertia for rotation about the x -axis, et cetera. Note that the previous expression consists of three terms that are quadratic in a momentum component. Hence, according to the law of equipartition of energy, the mean rotational energy of the molecule should be

$$\frac{1}{2} k_B T + \frac{1}{2} k_B T + \frac{1}{2} k_B T = \frac{3}{2} k_B T. \quad (5.192)$$

In fact, this is not the case. The reason for the discrepancy is that one of the principal axes of rotation of a hydrogen molecule corresponds to the axis that passes through the nuclei of the two

hydrogen atoms that constitute the molecule. The principal moment of inertia for rotation about this axis is much smaller than the principal moments of inertia for rotation about the other two principal axes. In fact, the former moment of inertia is of order $m_e d^2$, where m_e is the mass of an electron, and d the radius of a hydrogen atom (the contribution of the protons to the moment is negligible), whereas the latter two moments of inertia are of order $m_p D^2$, where m_p is the mass of a proton, and D the length of the atomic bond joining the two hydrogen atoms. Given that $m_e \ll m_p$, while $d \sim D$, the former moment of inertia is indeed much smaller than the latter two. It turns out that quantum mechanical considerations prevent a degree of rotational freedom with an anomalously small moment of inertia from contributing $(1/2)k_B T$ to the mean energy of the molecule. (See Section 5.5.8.) Hence, the mean rotational energy of a hydrogen molecule (or any diatomic molecule) is

$$\frac{1}{2}k_B T + \frac{1}{2}k_B T = k_B T. \quad (5.193)$$

According to the previous discussion, the mean energy of a hydrogen molecule is

$$\frac{3}{2}k_B T + k_B T = \frac{5}{2}k_B T, \quad (5.194)$$

where the former contribution is the molecule's mean translation kinetic energy, whereas the latter contribution is the molecule's mean rotational kinetic energy. Hence, the internal energy of a gas consisting of ν moles of hydrogen molecules is [see Equation (5.189)]

$$U = \nu N_A \frac{5}{2} k_B T = \nu \frac{5}{2} R T = \nu c_V T. \quad (5.195)$$

It follows that the molar specific heat capacity at constant volume of a hydrogen gas (or any diatomic gas) is

$$c_V = \frac{5}{2} R, \quad (5.196)$$

in accordance with Equation (5.107).

5.3.7 Brownian Motion

In 1827, Robert Brown was studying pollen grains of the plant *Clarkia pulchella* suspended in water under a microscope when he observed minute particles, ejected by the pollen grains, executing a jittery motion. Let us examine this phenomenon, which is known as *Brownian motion*.

Consider a particle of mass m that is suspended in a fluid. Let us investigate the motion of this particle parallel to the x -axis. The particle is subject to two types of force. First, a set of impulsive forces due to molecular bombardment. Second, a retarding force that is proportional to the particle's instantaneous speed through the surrounding fluid. Thus, the particle's equation of motion along the x -axis can be written

$$m \frac{d^2 x}{dt^2} = X(t) - \alpha \frac{dx}{dt}, \quad (5.197)$$

where $X(t)$ is the impulsive force due to molecular bombardment, and $-\alpha dx/dt$ the retarding force. It follows that

$$m x \frac{d^2 x}{dt^2} = x X - \alpha x \frac{dx}{dt}, \quad (5.198)$$

which can also be written

$$m \frac{d}{dt} \left(x \frac{dx}{dt} \right) - m \left(\frac{dx}{dt} \right)^2 = x X - \frac{\alpha}{2} \frac{dx^2}{dt}. \quad (5.199)$$

Taking the ensemble average of the previous equation, we obtain

$$m \frac{d}{dt} \left(\left\langle x \frac{dx}{dt} \right\rangle \right) - m \left\langle \left(\frac{dx}{dt} \right)^2 \right\rangle = \langle x X \rangle - \frac{\alpha}{2} \frac{d \langle x^2 \rangle}{dt}. \quad (5.200)$$

However,

$$\langle x X \rangle = 0, \quad (5.201)$$

because x and X are uncorrelated random variables whose mean values are zero. Furthermore,

$$\left\langle x \frac{dx}{dt} \right\rangle = 0, \quad (5.202)$$

because x and dx/dt are also uncorrelated random variables whose mean values are zero. Finally,

$$\frac{1}{2} m \left\langle \left(\frac{dx}{dt} \right)^2 \right\rangle = \frac{1}{2} k_B T, \quad (5.203)$$

by the law of equipartition of energy, where T is the temperature of the fluid. (See Section 5.3.4.) It follows from the previous four equations that

$$\frac{d \langle x^2 \rangle}{dt} = \frac{2 k_B T}{\alpha}, \quad (5.204)$$

which can be integrated to give

$$\langle x^2 \rangle = 2 D t, \quad (5.205)$$

where

$$D = \frac{k_B T}{\alpha}. \quad (5.206)$$

It can be seen, by comparison with the analysis of Sections 5.1.5 and 5.1.7, that molecular bombardment causes a particle immersed in a fluid to execute a random walk along the x -axis with diffusivity D .

Suppose that the particle is a sphere of radius a . Furthermore, suppose that the retarding force acting on the particle is due to fluid viscosity. According to *Stokes's law*,

$$\alpha = 6\pi \eta a, \quad (5.207)$$

where η is the viscosity of the fluid. It follows that

$$D = \frac{k_B T}{6\pi \eta a}. \quad (5.208)$$

This result, which was first obtained by Einstein in 1905, and was verified experimentally by Jean B. Perrin in 1910, served as the first convincing evidence of the existence of atoms and molecules. Note that the previous diffusivity scales as the inverse of the particle radius. Hence, only relatively small particles are likely to exhibit noticeable Brownian motion.

5.3.8 Mean Free Path

The *mean free path* is the average distance a molecule in a gas travels between collisions with other molecules. Let us crudely approximate the molecules in the gas as hard spheres of diameter R . Any two molecules whose centers are less than a distance R apart will collide. Suppose that one molecule is moving with velocity \mathbf{v} , whereas the other molecules are stationary. The moving molecule sweeps out a cylindrical volume $\pi R^2 \langle v \rangle$ in one second. Any other molecule whose center lies in this volume will collide with the moving molecule. There are $n \pi R^2 \langle v \rangle$ such molecules, where n is the number density of molecules. Hence, the number of collisions per second is

$$f = \pi R^2 n \langle v \rangle. \quad (5.209)$$

Thus, the mean distance that the molecule travels between collisions, which is the mean free path, is

$$l = \frac{\langle v \rangle}{f} = \frac{1}{\pi R^2 n}. \quad (5.210)$$

If we now take into account the fact that all of the molecules in the gas are moving then it is clear that the previous two equations generalize to give

$$f = \pi R^2 n \langle V \rangle, \quad (5.211)$$

$$l = \frac{\langle v \rangle}{\pi R^2 n \langle V \rangle}, \quad (5.212)$$

where \mathbf{V} is the relative velocity between molecules. Consider two molecules of velocities \mathbf{v}_1 and \mathbf{v}_2 . The relative velocity of the molecules is

$$\mathbf{V} = \mathbf{v}_1 - \mathbf{v}_2. \quad (5.213)$$

Now,

$$V^2 = v_1^2 + v_2^2 - 2 \mathbf{v}_1 \cdot \mathbf{v}_2, \quad (5.214)$$

which implies that

$$\langle V^2 \rangle = \langle v_1^2 \rangle + \langle v_2^2 \rangle - 2 \langle \mathbf{v}_1 \cdot \mathbf{v}_2 \rangle. \quad (5.215)$$

However, $\langle \mathbf{v}_1 \cdot \mathbf{v}_2 \rangle = 0$, because the cosine of the angle subtended between \mathbf{v}_1 and \mathbf{v}_2 is just as likely to be positive as to be negative. Thus, we deduce that

$$\langle V^2 \rangle = \langle v_1^2 \rangle + \langle v_2^2 \rangle = 2 \langle v^2 \rangle. \quad (5.216)$$

Assuming, as seems reasonable, that

$$\frac{\langle V \rangle}{\langle v \rangle} = \sqrt{\frac{\langle V^2 \rangle}{\langle v^2 \rangle}}, \quad (5.217)$$

we obtain

$$\langle V \rangle = \sqrt{2} \langle v \rangle. \quad (5.218)$$

Hence, Equation (5.212) yields

$$l = \frac{1}{\sqrt{2} \pi R^2 n}. \quad (5.219)$$

Let us estimate the mean free path in air at standard temperature ($T = 15^\circ \text{C}$) and pressure ($p = 10^5 \text{ N m}^{-2}$.) From the idea gas law, (5.97),

$$n = \frac{\nu N_A}{V} = \frac{p}{k_B T} = \frac{10^5}{(1.381 \times 10^{-23})(288)} = 2.5 \times 10^{25} \text{ m}^{-3}. \quad (5.220)$$

Now, $R = 2 \times 10^{-10} \text{ m}$ is a typical diameter of an air molecule. Thus, we obtain

$$l = \frac{1}{\sqrt{2} \pi (2 \times 10^{-10})^2 (2.5 \times 10^{25})} = 2 \times 10^{-7} \text{ m}. \quad (5.221)$$

Consider a molecule moving along the x -axis. The molecule is subject to random collisions. Thus, the probability that the molecule undergoes a collision between moving a distance x and moving a distance $x + dx$ is αdx , where α is a constant. Let $P(x)$ be the probability that the molecule moves a distance x without undergoing a collision. It is evident that the probability that the molecule's first collision occurs between moving a distance x and a distance $x + dx$ is $-(dP/dx) dx$. However, this probability is also equal to the probability that the molecule does not undergo a collision in moving a distance x , and then undergoes a collision between moving a distance x and a distance $x + dx$. In other words,

$$-\frac{dP}{dx} dx = [P(x)] [\alpha dx], \quad (5.222)$$

or

$$\frac{dP}{dx} = -\alpha P. \quad (5.223)$$

The previous equation can be integrated to give

$$P(x) = p_0 e^{-\alpha x}, \quad (5.224)$$

where p_0 is an arbitrary constant. However, $P(0) = 0$, because the molecule has no chance of undergoing a collision in moving zero distance. Hence, the probability that the molecule moves a distance x without undergoing a collision is

$$P(x) = e^{-\alpha x}. \quad (5.225)$$

Moreover, the probability that the molecule undergoes its first collision between moving a distance x and moving a distance $x + dx$ is $f(x) dx = [P(x)] [\alpha dx]$ (i.e., the molecule needs to not undergo a collision in moving a distance x , and then undergo a collision between moving a distance x and a distance $x + dx$), so

$$f(x) = \alpha e^{-\alpha x}. \quad (5.226)$$

The mean distance that the molecule travels before undergoing its first collision is (see Section 5.1.6)

$$\langle x \rangle = \int_0^{\infty} f(x) x dx = \int_0^{\infty} \alpha e^{-\alpha x} x dx = \frac{1}{\alpha} \int_0^{\infty} y e^{-y} dy = \frac{1}{\alpha}. \quad (5.227)$$

However, $\langle x \rangle$ is equivalent to the mean free path, l . Hence, the probability density for a molecule to move a distance x between collisions is

$$f(x) = \frac{e^{-x/l}}{l}. \quad (5.228)$$

Moreover, the probability that the molecule moves a distance x without undergoing a collision is

$$P(x) = e^{-x/l}. \quad (5.229)$$

5.3.9 Diffusion

Consider an ideal gas of uniform temperature, T , that has a number density gradient along the z -axis, such that

$$n(z) = n_0 + \frac{\partial n}{\partial z} dz. \quad (5.230)$$

Let $F(v)$ be the distribution of molecular speeds. Repeating the analysis of Section 5.3.2, the number of molecules per unit area, per second, whose speeds lie between v and $v + dv$, and whose directions of motion subtend an angle lying between θ and $\theta + d\theta$ with the z -axis, that cross the x - y plane is

$$dJ_z = [n' F(v) dv] [g(\theta) d\theta] [v_z] = [n' F(v) dv] \left[\frac{1}{2} \sin \theta d\theta \right] [v \cos \theta], \quad (5.231)$$

where n' is the number density where the molecules last made a collision. [See Equation (5.167).] On average, the molecules move a distance l (i.e., the mean free path) between collisions. Hence, $dz = -l \cos \theta$, and

$$n' = n_0 - \frac{\partial n}{\partial z} l \cos \theta. \quad (5.232)$$

Thus, the net flux of molecules across the x - y plane is

$$J_z = \frac{1}{2} \int_0^{\pi} \left[n_0 - \frac{\partial n}{\partial z} l \cos \theta \right] \cos \theta \sin \theta d\theta \int_0^{\infty} F(v) dv, \quad (5.233)$$

which gives

$$J_z = -D \frac{\partial n}{\partial z}, \quad (5.234)$$

where

$$D = \frac{1}{3} l \langle v \rangle. \quad (5.235)$$

Here, $\langle v \rangle$ is the mean molecular speed. (See Section 5.5.9.) Thus, we conclude that the flux of molecules in the z -direction is proportional to minus the local number density gradient along the z -axis. This result is known as *Fick's law*, after Adolf Fick who discovered it experimentally in 1855.

Consider a slab of gas lying between z and $z+dz$. The rate of change of the number of molecules contained in the slab is the difference between the flux of molecules into the slab and the flux of molecules out of the slab. In other words,

$$\frac{\partial(n A dz)}{\partial t} = [J_z(z, t) - J_z(z + dz, t)] A, \quad (5.236)$$

where A is the cross-sectional area of the slab. The previous equation implies that

$$\frac{\partial n}{\partial t} = -\frac{\partial J_z}{\partial z}. \quad (5.237)$$

However, making use of Fick's law, (5.234), we obtain

$$\frac{\partial n}{\partial t} = D \frac{\partial^2 n}{\partial z^2}. \quad (5.238)$$

The previous equation is known as the *diffusion equation*, and the constant D is known as the *diffusivity*.

It can be seen, by inspection, that one solution of the diffusion equation is

$$n(z, t) = n_0 + \frac{\delta n_0}{\sqrt{4\pi D t}} \exp\left(-\frac{z^2}{4 D t}\right), \quad (5.239)$$

where n_0 and δn_0 are arbitrary constants. Note that at $t = 0$,

$$n(z, 0) = n_0 + \delta n_0 \delta(z), \quad (5.240)$$

where $\delta(z)$ is a delta function. (See Section 2.1.6.) Moreover, at large times,

$$n(z, t \rightarrow \infty) = n_0. \quad (5.241)$$

Thus, our solution describes an initially localized Gaussian (see Section 5.1.7) density perturbation that gradually spreads out, and eventually disperses entirely. It is easily demonstrated that the width (i.e., standard deviation) of the density perturbation, σ_z , grows in time as

$$\sigma_z = \sqrt{2 D t}. \quad (5.242)$$

On the other hand, the maximum height of the perturbation decays in time as

$$h_z = \frac{\delta n_0}{\sqrt{4\pi D t}}. \quad (5.243)$$

Moreover, the area under the perturbation remains fixed as it evolves in time, which implies that the number of molecules associated with the density perturbation also remains fixed, as has to be the case (because we have not discussed any processes that create or destroy molecules). It is clear, from Sections 5.1.5 and 5.1.7, that the spreading of the density perturbation is due to a random walk of the excess molecules along the z -axis, under the action of molecular collisions.

Let us estimate the particle diffusivity in air at standard temperature ($T = 15^\circ \text{C}$) and pressure ($p = 10^5 \text{ N m}^{-2}$). The mean thermal speed of molecules of mass m in an ideal gas of temperature T is

$$\langle v \rangle = \sqrt{\frac{8}{\pi} \frac{k_B T}{m}}. \quad (5.244)$$

(See Section 5.5.9.) Hence, it follows from Equations (5.219), (5.220), and (5.235) that

$$D = \frac{2}{3\pi^{3/2}} \frac{1}{R^2 p} \sqrt{\frac{(k_B T)^3}{m}}, \quad (5.245)$$

where R is the molecular diameter. Thus, the diffusivity scales as $1/p$ at constant temperature, as $T^{3/2}$ at constant pressure, and as $T^{1/2}$ at constant volume. Given that $m = 29 m_p$ for air, where m_p is the proton mass, and $R = 2 \times 10^{-10} \text{ m}$, we deduce that

$$D = 3 \times 10^{-5} \text{ m}^2 \text{ s}^{-1}. \quad (5.246)$$

This is a very small diffusivity. According to Equation (5.242), it takes about 4.6 hours for a molecule to diffuse a distance of 1 meter in air.

5.3.10 Thermal Conductivity

Consider a gas of uniform number density, n , that has a temperature gradient along the z -axis, such that

$$T(z) = T_0 + \frac{\partial T}{\partial z} dz. \quad (5.247)$$

Let each molecule in the gas have a mean thermal energy $\epsilon(T)$ (under most circumstances this energy is the sum of the molecule's translational and rotational kinetic energy). Slightly modifying the analysis of the previous section, the thermal energy per unit area, per second, carried by molecules whose speeds lie between v and $v + dv$, and whose directions of motion subtend an angle lying between θ and $\theta + d\theta$ with the z -axis, that cross the x - y plane, is

$$dq_z = [\epsilon'] [n F(v) dv] \left[\frac{1}{2} \sin \theta d\theta \right] [v \cos \theta], \quad (5.248)$$

where $F(v)$ is the distribution of molecular speeds, and $\epsilon' = \epsilon(T')$, where T' is the temperature where the molecules last made a collision. On average, the molecules move a distance l (i.e., the mean free path) between collisions. Hence, $dz = -l \cos \theta$, and

$$T' = T_0 - \frac{\partial T}{\partial z} l \cos \theta. \quad (5.249)$$

which implies that

$$\epsilon' \simeq \epsilon_0 - C \frac{\partial T}{\partial z} l \cos \theta, \quad (5.250)$$

where

$$C = \frac{\partial \epsilon}{\partial T} \quad (5.251)$$

is the specific heat per molecule. Thus, the net heat flux across the x - y plane is

$$q_z = \frac{n}{2} \int_0^\pi \left[\epsilon_0 - C \frac{\partial T}{\partial z} l \cos \theta \right] \cos \theta \sin \theta d\theta \int_0^\infty F(v) dv, \quad (5.252)$$

which gives

$$q_z = -\kappa \frac{\partial T}{\partial z}, \quad (5.253)$$

where

$$\kappa = \frac{1}{3} n l C \langle v \rangle \quad (5.254)$$

is termed the *thermal conductivity*. Here, $\langle v \rangle$ is the mean molecular speed. (See Section 5.5.9.) Thus, we conclude that the heat flux in the z -direction is proportional to minus the local temperature gradient along the z -axis. This is another example of Fick's law.

Equations (5.219), (5.244), and (5.253) yield

$$\kappa = \frac{2}{3} \frac{C}{\pi^{3/2}} \frac{1}{R^2} \sqrt{\frac{k_B T}{m}}, \quad (5.255)$$

where R is the molecular diameter, and m the molecular mass. Moreover, for a diatomic gas, such as air, $C = (5/2) k_B$. (See Section 5.3.6.) It can be seen that the thermal conductivity of an ideal gas scales as $T^{1/2}$, and is independent of the pressure at constant temperature. The previous formula yields the following estimate for the thermal conductivity of air (assuming that $T = 15^\circ \text{C}$, $R = 2 \times 10^{-10} \text{m}$, and $m = 29 m_p$),

$$\kappa = 3 \times 10^{-2} \text{W m}^{-1} \text{K}^{-1}, \quad (5.256)$$

which turns out to be fairly accurate

Consider a slab of gas lying between z and $z + dz$. The rate of change of the thermal energy contained in the slab is the difference between the flux of heat into the slab and the flux of heat out of the slab. In other words,

$$\frac{\partial(n A dz C T)}{\partial t} = [q_z(z, t) - q_z(z + dz, t)] A, \quad (5.257)$$

where A is the cross-sectional area of the slab. It follows that

$$\frac{\partial T}{\partial t} = -\frac{1}{nC} \frac{\partial q_z}{\partial z}. \quad (5.258)$$

Making use of Equations (5.253) and (5.254), we obtain

$$\frac{\partial T}{\partial t} = D_\kappa \frac{\partial^2 T}{\partial z^2}, \quad (5.259)$$

where

$$D_\kappa = \frac{\kappa}{nC} = \frac{1}{3} l \langle v \rangle. \quad (5.260)$$

Of course, Equation (5.259) is the diffusion equation, and D_κ is the associated diffusivity. Thus, we conclude, by comparison with the analysis in the previous section, that heat diffuses through an ideal gas at the same very slow rate at which molecules diffuse. Moreover, it is clear, from Sections 5.1.5 and 5.1.7, that heat diffusion is due to a random walk of molecules with excess energy under the action of molecular collisions.

5.3.11 Viscosity

Consider a gas of uniform number density, n , that has a gradient in the x -component of its mean flow velocity, V_x , along the z -axis, such that

$$V_x(z) = V_{x0} + \frac{\partial V_x}{\partial z} dz. \quad (5.261)$$

It is assumed that the mean flow velocity is much less than the mean molecular speed. Slightly modifying the analysis of the previous two sections, the x -momentum per unit area, per second, carried by molecules whose speeds lie between v and $v+dv$, and whose directions of motion subtend an angle lying between θ and $\theta + d\theta$ with the z -axis, that cross the x - y plane, is

$$dP_{xz} = [m V'_x] [n F(v) dv] \left[\frac{1}{2} \sin \theta d\theta \right] [v \cos \theta], \quad (5.262)$$

where $F(v)$ is the distribution of molecular speeds, m the molecular mass, and V'_x the x -component of flow velocity where the molecules last made a collision. On average, the molecules move a distance l (i.e., the mean free path) between collisions. Hence, $dz = -l \cos \theta$, and

$$V'_x = V_{x0} - \frac{\partial V_x}{\partial z} l \cos \theta. \quad (5.263)$$

Thus, the net flux of x -momentum across the x - y plane is

$$P_{xz} = \frac{mn}{2} \int_0^\pi \left[V_{x0} - \frac{\partial V_x}{\partial z} l \cos \theta \right] \cos \theta \sin \theta d\theta \int_0^\infty F(v) dv, \quad (5.264)$$

which gives

$$P_{xz} = -\eta \frac{\partial V_x}{\partial z}, \quad (5.265)$$

where

$$\eta = \frac{1}{3} mn l \langle v \rangle \quad (5.266)$$

is termed the *viscosity*. Here, $\langle v \rangle$ is the mean molecular speed. (See Section 5.5.9.) Thus, we conclude that the flux of x -momentum in the z -direction is proportional to minus the gradient of the x -component of the flow velocity with respect to z . This is yet another example of Fick's law.

Equations (5.219), (5.244), and (5.266) imply that yield

$$\eta = \frac{2}{3} \frac{1}{\pi^{3/2}} \frac{1}{R^2} \sqrt{k_B T m}, \quad (5.267)$$

where R is the molecular diameter, and T the gas temperature. It can be seen that the viscosity of an ideal gas scales as $T^{1/2}$, and is independent of the pressure at constant temperature. The previous formula yields the following estimate for the viscosity of air (assuming that $T = 15^\circ \text{C}$, $R = 2 \times 10^{-10} \text{m}$, and $m = 29 m_p$),

$$\eta = 4 \times 10^{-5} \text{N s m}^{-2}, \quad (5.268)$$

which turns out to be too large by a factor 2 (because of the approximate nature of our calculation).

Consider a slab of gas lying between z and $z + dz$. The rate of change of the x -momentum contained in the slab is the difference between the flux of momentum into the slab and the flux of momentum out of the slab. In other words,

$$\frac{\partial(n A dz m V_x)}{\partial t} = [P_{xz}(z, t) - P_{xz}(z + dz, t)] A, \quad (5.269)$$

where A is the cross-sectional area of the slab. It follows that

$$\frac{\partial V_x}{\partial t} = -\frac{1}{n m} \frac{\partial P_{xz}}{\partial z}. \quad (5.270)$$

Making use of Equations (5.265) and (5.266), we get

$$\frac{\partial V_x}{\partial t} = D_\eta \frac{\partial^2 V_x}{\partial z^2}, \quad (5.271)$$

where

$$D_\eta = \frac{\eta}{n m} = \frac{1}{3} l \langle v \rangle. \quad (5.272)$$

Of course, Equation (5.271) is the diffusion equation, and D_η is the associated diffusivity. Thus, we conclude, by comparison with the analysis in the previous two sections, that momentum diffuses through an ideal gas at the same very slow rate at which molecules and heat diffuse. Moreover, it is clear, from Sections 5.1.5 and 5.1.7, that momentum diffusion is due to a random walk of molecules with excess momentum under the action of molecular collisions.

Equations (5.235), (5.260), and (5.272) lead to the prediction that

$$D = \frac{\kappa}{n C} = \frac{\eta}{n m} \quad (5.273)$$

for an ideal gas, where D is the molecular diffusivity, κ the thermal conductivity, η the viscosity, n the number density of molecules, m the molecular mass, and C the molecular heat capacity. It turns

out that this prediction is only approximately true due to additional factors, such as intermolecular forces, that have not been incorporated into our highly simplified analysis. For example, the ratio $\kappa m/(\eta C)$, which should be unity according to simple kinetic theory, actually takes the values 2.40, 2.49, 1.91, 1.91, and 1.90 for helium, argon, hydrogen, nitrogen, and oxygen, respectively, at standard temperature and pressure.

The prediction that the thermal conductivity and viscosity of an ideal gas are both independent of the gas pressure, at constant temperature, breaks down when the pressure becomes sufficiently low that the mean free path between collisions becomes comparable with the size of the gas's container. Under these circumstances, the thermal conductivity and viscosity both become approximately proportional to the pressure, at constant temperature.

5.3.12 Molecular Flow

Consider the flow of an ideal gas down a uniform pipe of circular cross-section in the limit that the gas pressure is sufficiently low that the mean free path between collisions greatly exceeds the diameter of the pipe. This type of flow is known as *molecular flow*.

Suppose that the pipe runs along the z -axis, and that a pressure difference, Δp , is established between the two ends of the pipe, in order to drive the flow. The temperature of the gas is assumed to be uniform along the pipe. Finally, the length of the pipe, L , is assumed to be much greater than its diameter, d . Making use of Equation (5.170), the net flux of molecules down the pipe at position z is

$$\Phi_z(z) = \frac{1}{4} \langle v \rangle [n(z-d) - n(z+d)], \quad (5.274)$$

where $\langle v \rangle$ is the mean molecular speed [which is constant because it only depends on the temperature; see Equation (5.433)], and $n(z)$ the molecular number density. The right-hand side of the previous equation represents the difference between the particle flux in the $+z$ -direction and that in the $-z$ -direction. The former flux is characterized by the value the number density calculated at the position at which the molecules moving in the $+z$ direction last collided with the wall of the pipe, which is estimated to be $z-d$. Likewise, the latter flux is characterized by the number density calculated at the position at which the molecules moving in the $-z$ direction last collided with the wall of the pipe, which is estimated to be $z+d$. According to Equation (5.175), the pressure of the gas in the pipe is

$$p(z) = \frac{1}{3} n(z) m \langle v^2 \rangle, \quad (5.275)$$

where m is the molecular mass, and $\langle v^2 \rangle$ is the mean square molecular speed (which is also constant because it only depends on the temperature). [See Equation (5.434).] The previous two equations yield

$$\begin{aligned} \Phi_z(z) &= \frac{3}{4} \frac{\langle v \rangle}{m \langle v^2 \rangle} [p(z-d) - p(z+d)] \\ &\simeq -\frac{3}{2} \frac{d \langle v \rangle}{m \langle v^2 \rangle} \frac{dp}{dz}. \end{aligned} \quad (5.276)$$

In a steady state, Φ_z must be uniform along the pipe, which implies that dp/dz is also uniform. Hence, we can write

$$-\frac{dp}{dz} = \frac{\Delta p}{L}, \quad (5.277)$$

which yields

$$\Phi_z = \frac{3}{2} \frac{\langle v \rangle}{\langle v^2 \rangle} \frac{d}{mL} \Delta p. \quad (5.278)$$

Now, the cross-sectional area of the pipe is $A = \pi d^2/4$. Hence, the rate of mass flow down the pipe is

$$\dot{M} = \Phi_z A m = \frac{3\pi}{8} \frac{\langle v \rangle}{\langle v^2 \rangle} \frac{d^3}{L} \Delta p. \quad (5.279)$$

However, in an ideal gas,

$$\langle v \rangle^2 = \frac{8}{3\pi} \langle v^2 \rangle. \quad (5.280)$$

(See Section 5.5.9.) Thus, we obtain

$$\dot{M} = \frac{d^3}{\langle v \rangle L} \Delta p. \quad (5.281)$$

Given that $\langle v \rangle \propto T^{1/2}$ [see Equation (5.244)], where T is the temperature of the gas, we deduce that the mass flow rate due to molecular flow of an ideal gas down a pipe, when a given pressure difference is established between the two ends of the pipe, is proportional to the cube of the pipe diameter, inversely proportional to the length of the pipe, and inversely proportional to the square-root of the temperature.

Now, the standard formula for the mass flow rate of a viscous fluid down a pipe of circular cross-section, which applies to the case under discussion when the mean free path between collisions is much less than the diameter of the pipe, is

$$\dot{M} = \frac{\pi}{128} \frac{d^4 n m}{\eta L} \Delta p. \quad (5.282)$$

Making use of Equation (5.266), we deduce that the mass flow rate of an ideal gas down a pipe of circular cross-section, in the short mean free path limit, is

$$\dot{M} = \frac{3\pi}{128} \frac{d^4}{\langle v \rangle l L} \Delta p, \quad (5.283)$$

where l is the mean free path. Thus, Equation (5.281) holds when $l \gg d$, and Equation (5.283) holds when $l \ll d$.

5.3.13 Molecular Effusion

Consider an ideal gas held in a container that is divided in two by a partition. Let the gas to the left of the partition have temperature T_1 , pressure p_1 , and number density n_1 . Likewise, let the gas to the right of the partition have temperature T_2 , pressure p_2 , and number density n_2 . Suppose that the

partition contains a small hole of cross-sectional area A . Suppose that the dimension of the hole is much less than the mean free path between collisions for the gases on either side of the partition. In this case, as described in Section 5.3.2, the gases effuse through the hole, rather than flowing through it in a hydrodynamical fashion.

According to Equations (5.172), (5.244), and the ideal gas law, which states that $p_1 = n_1 k_B T_1$, the number of molecules per unit time that effuse from the left partition to the right partition is

$$\dot{N}_{12} = \frac{A}{\sqrt{2\pi m}} \frac{p_1}{\sqrt{k_B T_1}}, \quad (5.284)$$

where m is the molecular mass. Likewise, the number of molecules per unit time that effuse from the right partition to the left partition is

$$\dot{N}_{21} = \frac{A}{\sqrt{2\pi m}} \frac{p_2}{\sqrt{k_B T_2}}. \quad (5.285)$$

In a steady state, we require $\dot{N}_{12} = \dot{N}_{21}$. Hence, we deduce that

$$\frac{p_1}{p_2} = \left(\frac{T_1}{T_2} \right)^{1/2}. \quad (5.286)$$

Thus, in equilibrium, a higher pressure prevails in the part of the container held at a higher temperature. This result is different to that we would obtain in limit in which the mean free path is much less than the size of the hole. In this limit, the gases on either side of the hole flow through it in a hydrodynamical fashion, and an equilibrium state is achieved when

$$p_1 = p_2. \quad (5.287)$$

5.4 Statistical Mechanics

5.4.1 Specification of State of Many-Particle System

Let us consider how we might specify the state of a system consisting of a great many particles, such as an ideal gas. Consider the simplest possible dynamical system, which consists of a single spinless particle moving classically in one dimension. Assuming that we know the particle's equation of motion, the state of the system is fully specified once we simultaneously measure the particle's displacement, q , and momentum, p . In fact, if we know q and p then we can calculate the state of the system at all subsequent times using the equation of motion. In practice, it is impossible to specify q and p exactly, because there is always an intrinsic error in any experimental measurement.

Consider the time evolution of q and p . This can be visualized by plotting the point (q, p) in the q - p plane. This plane is generally known as *phase-space*. In general, as time progresses, the point (q, p) will trace out some very complicated pattern in phase-space. Suppose that we divide phase-space into rectangular cells of uniform dimensions δq and δp . Here, δq is the intrinsic error

in the position measurement, and δp the intrinsic error in the momentum measurement. The area of each cell is

$$\delta q \delta p = h_0, \quad (5.288)$$

where h_0 is a small constant having the dimensions of angular momentum. The coordinates q and p can now be conveniently specified by indicating the cell in phase-space into which they plot at any given time. This procedure automatically ensures that we do not attempt to specify q and p to an accuracy greater than our experimental error, which would clearly be pointless.

Let us now consider a single spinless particle moving in three dimensions. In order to specify the state of the system, we now need to know three q - p pairs; that is, q_x - p_x , q_y - p_y , and q_z - p_z . Incidentally, the number of q - p pairs needed to specify the state of the system is usually called the *number of degrees of freedom* of the system. Thus, a single particle moving in one dimension constitutes a one degree of freedom system, whereas a single particle moving in three dimensions constitutes a three degree of freedom system.

Consider the time evolution of \mathbf{q} and \mathbf{p} , where $\mathbf{q} = (q_x, q_y, q_z)$, et cetera. This can be visualized by plotting the point (\mathbf{q}, \mathbf{p}) in the six-dimensional \mathbf{q} - \mathbf{p} phase-space. Suppose that we divide the q_x - p_x plane into rectangular cells of uniform dimensions δq and δp , and do likewise for the q_y - p_y and q_z - p_z planes. Here, δq and δp are again the intrinsic errors in our measurements of position and momentum, respectively. This is equivalent to dividing phase-space up into regular six-dimensional cells of volume h_0^3 . The coordinates \mathbf{q} and \mathbf{p} can now be conveniently specified by indicating the cell in phase-space into which they plot at any given time. Again, this procedure automatically ensures that we do not attempt to specify \mathbf{q} and \mathbf{p} to an accuracy greater than our experimental error.

Finally, let us consider a system consisting of N spinless particles moving classically in three dimensions. In order to specify the state of the system, we need to specify a large number of q - p pairs. The requisite number is simply the number of degrees of freedom, f . For the present case, $f = 3N$, because each particle needs three q - p pairs. Thus, phase-space (i.e., the space of all the q - p pairs) now possesses $2f = 6N$ dimensions. Consider a particular pair of phase-space coordinates, q_i and p_i . As before, we divide the q_i - p_i plane into rectangular cells of uniform dimensions δq and δp . This is equivalent to dividing phase-space into regular $2f$ dimensional cells of volume h_0^f . The state of the system is specified by indicating which cell it occupies in phase-space at any given time.

In principle, we can specify the state of the system to arbitrary accuracy, by taking the limit $h_0 \rightarrow 0$. In reality, we know from Heisenberg's uncertainty principle (see Section 4.2.7) that it is impossible to simultaneously measure a coordinate, q_i , and its associated momentum, p_i , to greater accuracy than $\delta q_i \delta p_i = \hbar/2$. Here, \hbar is Planck's constant divided by 2π . This implies that

$$h_0 \geq \hbar/2. \quad (5.289)$$

In other words, the uncertainty principle sets a lower limit on how finely we can chop up classical phase-space.

In quantum mechanics, we can specify the state of the system by giving its wavefunction at time t ,

$$\psi(q_1, \dots, q_f, s_1, \dots, s_g, t), \quad (5.290)$$

where f is the number of translational degrees of freedom, and g the number of internal (e.g., spin) degrees of freedom. For instance, if the system consists of N spin-one-half particles then there will be $3N$ translational degrees of freedom, and N spin degrees of freedom (because the spin of each particle can either be directed up or down along the z -axis). Alternatively, if the system is in a *stationary state* (see Section 4.2.9) then we can just specify $f + g$ quantum numbers. Either way, the future time evolution of the wavefunction is fully determined by Schrödinger's equation. In reality, this approach is not practical because Schrödinger's equation for the system is only known approximately. Typically, we are dealing with a system consisting of many weakly-interacting particles. We usually know Schrödinger's equation for completely non-interacting particles, but the component of the equation associated with particle interactions is either impossibly complicated, or not very well known. We can define approximate stationary eigenstates using the Schrödinger's equation for non-interacting particles. The state of the system is then specified by the quantum numbers identifying these eigenstates. In the absence of particle interactions, if the system starts off in a stationary state then it stays in that state for ever, so its quantum numbers never change. The interactions allow the system to make transitions between different "stationary" states, causing its quantum numbers to change in time.

5.4.2 Principle of Equal A Priori Probabilities

We now know how to specify the instantaneous state of a many-particle system. In principle, such a system is completely deterministic. If we know the initial state, and the equations of motion, then we can evolve the system forward in time, and, thereby, determine all future states. In reality, it is quite impossible to specify the initial state, or the equations of motion, to sufficient accuracy for this method to have any chance of working. Furthermore, even if it were possible, it would still not be a practical proposition to evolve the equations of motion. We are typically dealing with systems containing Avogadro's number of particles; that is, about 10^{24} particles. We cannot evolve 10^{24} simultaneous differential equations. Even if we could, we would not want to. After all, we are not particularly interested in the motions of individual particles. What we really require is statistical information regarding the motions of all particles in the system.

Clearly, what is needed here is a statistical treatment of the problem. Instead of focusing on a single system, let us proceed, in the usual manner, and consider a statistical ensemble consisting of a large number of identical systems. (See Section 5.1.1.) In general, these systems are distributed over many different states at any given time. In order to evaluate the probability that the system possesses a particular property, we merely need to find the number of systems in the ensemble that exhibit this property, and then divide by the total number of systems, in the limit as the latter number tends to infinity.

We can usually place some general constraints on the system. Typically, we know the total internal energy, U , the total volume, V , and the total number of particles, N . To be more exact, we can only really say that the total internal energy lies between U and $U + \delta U$, et cetera, where δU is an experimental error. Thus, we need only concern ourselves with those systems in the ensemble exhibiting states that are consistent with the known constraints. We shall call these the *states accessible to the system*. In general, there are a great many such states.

We now need to calculate the probability of the system being found in each of its accessible

states. In fact, the only way that we could calculate these probabilities would be to evolve all of the systems in the ensemble in time, and observe how long, on average, they spend in each accessible state. But, as we have already discussed, such a calculation is completely out of the question. Instead, we shall effectively guess the probabilities.

Let us consider an isolated system in equilibrium. In this situation, we would expect the probability of the system being found in one of its accessible states to be independent of time. This implies that the statistical ensemble does not evolve with time. Individual systems in the ensemble will constantly change state, but the average number of systems in any given state should remain constant. Thus, all macroscopic parameters describing the system, such as the internal energy and the volume, should also remain constant. There is nothing in the laws of mechanics that would lead us to suppose that the system will be found more often in one of its accessible states than in another. We assume, therefore, that *the system is equally likely to be found in any of its accessible states*. This assumption is called the *principle of equal a priori probabilities*, and lies at the heart of statistical mechanics. In fact, we use assumptions like this all of the time without really thinking about them. Suppose that we were asked to pick a card at random from a well-shuffled pack of ordinary playing cards. Most people would accept that we have an equal probability of picking any card in the pack. There is nothing that would favor one particular card over all of the others. Hence, because there are fifty-two cards in a normal pack, we would expect the probability of picking the ace of spades, say, to be $1/52$. We could now place some constraints on the system. For instance, we could only count red cards, in which case the probability of picking the ace of hearts, say, would be $1/26$, by the same reasoning. In both cases, we have used the principle of equal a priori probabilities. In statistical mechanics, we treat a many-particle system a little like an extremely large pack of cards. Each accessible state corresponds to one of the cards in the pack. The interactions between particles cause the system to continually change state. This is equivalent to constantly shuffling the pack. Finally, an observation of the state of the system is like picking a card at random from the pack. The principle of equal a priori probabilities then boils down to saying that we have an equal chance of choosing any particular card.

5.4.3 Probability Calculations

The principle of equal a priori probabilities is fundamental to all of statistical mechanics, and allows a complete description of the properties of macroscopic systems in equilibrium. Consider a system in equilibrium that is isolated, so that its total internal energy is known to have a constant value lying somewhere in the range U to $U + \delta U$. In order to make statistical predictions, we focus attention on an ensemble of such systems, all of which have their internal energy in this range. Let $\Omega(U)$ be the total number of different states in the ensemble with internal energies in the specified range. Suppose that, among these states, there are a number $\Omega(U; x_k)$ for which some parameter, x , of the system assumes the discrete value x_k . (This discussion can easily be generalized to deal with a parameter that can assume a continuous range of values.) The principle of equal a priori probabilities tells us that all of the $\Omega(U)$ accessible states of the system are equally likely to occur in the ensemble. It follows that the probability, $P(x_k)$, that the parameter x of the system assumes

the value x_k is simply

$$P(x_k) = \frac{\Omega(U; x_k)}{\Omega(U)}. \quad (5.291)$$

Clearly, the mean value of x for the system is given by

$$\langle x \rangle = \frac{\sum_k \Omega(U; x_k) x_k}{\Omega(U)}, \quad (5.292)$$

where the sum is over all possible values that x can assume.

5.4.4 Number of Accessible States of Ideal Gas

Consider an ideal gas, made up of spinless monatomic particles, whose volume is V , and whose internal energy lies in the range U to $U + \delta U$. Let $\Omega(U, V)$ be the total number of microscopic states that satisfy these constraints. This is a particularly simple example, because, for such a gas, the particles possess translational, but no internal (e.g., vibrational, rotational, or spin), degrees of freedom. By definition, interatomic forces are negligible in an ideal gas. In other words, the individual particles move in an approximately uniform potential. It follows that the internal energy of the gas is just the total translational kinetic energy of its constituent particles, so that

$$U = \frac{1}{2m} \sum_{i=1, N} \mathbf{p}_i^2, \quad (5.293)$$

where m is the particle mass, N the total number of particles, and \mathbf{p}_i the vector momentum of the i th particle.

Consider the system in the limit in which the internal energy, U , of the gas is much greater than the ground-state energy, so that all of the quantum numbers are large. The classical version of statistical mechanics, in which we divide up phase-space into cells of equal volume, is valid in this limit. (See Section 5.4.1.) The number of states, $\Omega(U, V)$, lying between the internal energies U and $U + \delta U$ is simply equal to the number of cells in phase-space contained between these energies. In other words, $\Omega(U, V)$ is proportional to the volume of phase-space between these two energies:

$$\Omega(U, V) \propto \int_U^{U+\delta U} d^3\mathbf{r}_1 \cdots d^3\mathbf{r}_N d^3\mathbf{p}_1 \cdots d^3\mathbf{p}_N. \quad (5.294)$$

Here, the integrand is the element of volume of phase-space, with

$$d^3\mathbf{r}_i \equiv dx_i dy_i dz_i, \quad (5.295)$$

$$d^3\mathbf{p}_i \equiv dp_{xi} dp_{yi} dp_{zi}, \quad (5.296)$$

where (x_i, y_i, z_i) and (p_{xi}, p_{yi}, p_{zi}) are the Cartesian coordinates and momentum components of the i th particle, respectively. The integration is over all coordinates and momenta such that the total internal energy of the system lies between U and $U + \delta U$.

For an ideal gas, the total internal energy U does not depend on the positions of the particles. [See Equation (5.293).] This implies that the integration over the position vectors, \mathbf{r}_i , can be

performed immediately. Because each integral over \mathbf{r}_i extends over the volume of the container (the particles are, of course, not allowed to stray outside the container), $\int d^3\mathbf{r}_i = V$. There are N such integrals, so Equation (5.294) reduces to

$$\Omega(U, V) \propto V^N \chi(U), \quad (5.297)$$

where

$$\chi(U) \propto \int_U^{U+\delta U} d^3\mathbf{p}_1 \cdots d^3\mathbf{p}_N \quad (5.298)$$

is a momentum-space integral that is independent of the volume.

The internal energy of the system can be written

$$U = \frac{1}{2m} \sum_{i=1, N} \sum_{\alpha=1, 3} p_{\alpha i}^2 \quad (5.299)$$

because $\mathbf{p}_i^2 = p_{1i}^2 + p_{2i}^2 + p_{3i}^2$, denoting the (x, y, z) components by $(1, 2, 3)$, respectively. The previous sum contains $3N$ square terms. For $U = \text{constant}$, Equation (5.299) describes the locus of a sphere of radius $R(U) = (2mU)^{1/2}$ in the $3N$ -dimensional space of the momentum components. Hence, $\chi(U)$ is proportional to the volume of momentum phase-space contained in the region lying between the sphere of radius $R(U)$, and that of slightly larger radius $R(U + \delta U)$. This volume is proportional to the area of the inner sphere multiplied by $\delta R \equiv R(U + \delta U) - R(U)$. Because the area varies like R^{3N-1} , and $\delta R \propto \delta U/U^{1/2}$, we have

$$\chi(U) \propto R^{3N-1}/U^{1/2} \propto U^{3N/2-1}. \quad (5.300)$$

Combining this result with Equation (5.297), we obtain

$$\Omega(U, V) \simeq B V^N U^{3N/2}, \quad (5.301)$$

where B is a constant independent of V or U , and we have also made use of the fact that $N \gg 1$ for a typical ideal gas. Note that $\Omega(U, V)$ is a very strongly increasing function of U because $N \gg 1$.

5.4.5 Thermal Interaction

Consider a purely thermal interaction between two systems, A and A' , that contain a large number of particles. Suppose that the internal energies of these two systems are U and U' , respectively. The external parameters are held fixed, so that systems A and A' cannot do work on one another. However, we shall assume that the systems are free to exchange heat energy (i.e., they are in thermal contact). It is convenient to divide the energy scale into small subdivisions of width δU . The number of accessible states of A (i.e., states in which the internal energy of the whole system lies between U and $U + \delta U$) is denoted $\Omega(U)$. Likewise, the number of accessible states of A' is denoted $\Omega'(U')$.

The combined system $A^{(0)} = A + A'$ is assumed to be isolated (i.e., it neither does work on, nor exchanges heat with, its surroundings). It follows the total internal energy, $U^{(0)}$, is constant.

When speaking of thermal contact between two distinct systems, we usually assume that the mutual interaction is sufficiently weak for the internal energies to be additive. Thus,

$$U + U' \simeq U^{(0)} = \text{constant}. \quad (5.302)$$

According to Equation (5.302), if the internal energy of A lies in the range U to $U + \delta U$ then the internal energy of A' must lie between $U^{(0)} - U - \delta U$ and $U^{(0)} - U$. Thus, the number of accessible states for each system is given by $\Omega(U)$ and $\Omega'(U^{(0)} - U)$, respectively. Because every possible state of A can be combined with every possible state of A' to form a distinct state, the total number of distinct states accessible to $A^{(0)}$ when the energy of A lies in the range U to $U + \delta U$ is

$$\Omega^{(0)}(U) = \Omega(U) \Omega'(U^{(0)} - U). \quad (5.303)$$

Consider an ensemble of pairs of thermally interacting systems, A and A' , that are left undisturbed, so that they can attain thermal equilibrium. The principle of equal a priori probabilities is applicable to this situation. (See Section 5.4.2.) According to this principle, the probability of occurrence of a given macroscopic state is proportional to the number of accessible microscopic states, because all microscopic states are equally likely. Thus, the probability that the system A has an energy lying in the range U to $U + \delta U$ can be written

$$P(U) = C \Omega(U) \Omega'(U^{(0)} - U), \quad (5.304)$$

where C is a constant that is independent of U .

We know, from Section 5.4.4, that the typical variation of the number of accessible states with energy is of the form

$$\Omega \propto U^{3N/2}, \quad (5.305)$$

where N is the number of molecules. For a macroscopic system, N is an exceedingly large number. It follows that the probability, $P(U)$, in Equation (5.304) is the product of an extremely rapidly increasing function of U , and an extremely rapidly decreasing function of U . Hence, we would expect the probability to exhibit a very pronounced maximum at some particular value of the energy, U_f .

5.4.6 Thermodynamic Temperature

Suppose that the systems A and A' are initially thermally isolated from one another, with respective internal energies U_i and U'_i . If the two systems are subsequently placed in thermal contact, so that they are free to exchange heat energy, then, in general, the resulting state is an extremely improbable one [i.e., $P(U_i)$ is much less than the peak probability]. The configuration will, therefore, tend to evolve in time until the two systems attain final energies, U_f and U'_f , which are such that $P(U_f)$ is maximized. In the special case where the initial energies, U_i and U'_i , lie very close to the final energies, U_f and U'_f , respectively, there is no change in the two systems when they are brought into thermal contact, because the initial state already corresponds to a state of maximum probability.

It follows from energy conservation that

$$U_f + U'_f = U_i + U'_i. \quad (5.306)$$

The energy change in each system is simply the net heat absorbed, so that

$$Q = U_f - U_i, \quad (5.307)$$

$$Q' = U'_f - U'_i. \quad (5.308)$$

The conservation of energy then reduces to

$$Q + Q' = 0. \quad (5.309)$$

In other words, the heat given off by one system is equal to the heat absorbed by the other. (In our notation, absorbed heat is positive, and emitted heat is negative.)

It is clear that if the systems A and A' are suddenly brought into thermal contact then they will only exchange heat, and evolve towards a new equilibrium state, if the final state is more probable than the initial one. In other words, the system will evolve if

$$P(U_f) > P(U_i), \quad (5.310)$$

or

$$\ln P(U_f) > \ln P(U_i), \quad (5.311)$$

because the logarithm is a monotonic function. The previous inequality can be written

$$\ln \Omega(U_f) + \ln \Omega'(U'_f) > \ln \Omega(U_i) + \ln \Omega'(U'_i), \quad (5.312)$$

with the aid of Equation (5.304). Taylor expansion to first order yields

$$\frac{\partial \ln \Omega(U_i)}{\partial U} (U_f - U_i) + \frac{\partial \ln \Omega'(U'_i)}{\partial U'} (U'_f - U'_i) > 0, \quad (5.313)$$

which finally gives

$$[\beta(U_i) - \beta'(U'_i)] Q > 0, \quad (5.314)$$

where

$$\beta = \frac{\ln \Omega}{\partial U}, \quad (5.315)$$

$$\beta' = \frac{\ln \Omega'}{\partial U'}, \quad (5.316)$$

and use has been made of Equations (5.307)–(5.309).

It is clear, from the previous analysis, that the parameter β has the following properties:

1. If two systems separately in equilibrium have the same value of β then the systems will remain in equilibrium when brought into thermal contact with one another.
2. If two systems separately in equilibrium have different values of β then the systems will not remain in equilibrium when brought into thermal contact with one another. Instead, the system with the higher value of β will absorb heat from the other system until the two β values are the same. [See Equation (5.314).]

Let us define the dimensionless parameter, T , such that

$$\frac{1}{k_B T} \equiv \beta = \frac{\partial \ln \Omega}{\partial E}, \quad (5.317)$$

where k_B is the Boltzmann constant. The parameter T is termed the *thermodynamic temperature*, and controls heat flow in much the same manner as a conventional temperature. Thus, if two isolated systems in equilibrium possess the same thermodynamic temperature then they will remain in equilibrium when brought into thermal contact. However, if the two systems have different thermodynamic temperatures then heat will flow from the system with the higher temperature (i.e., the “hotter” system) to the system with the lower temperature, until the temperatures of the two systems are the same. In addition, suppose that we have three systems, A , B , and C . We know that if A and B remain in equilibrium when brought into thermal contact then their temperatures are the same, so that $T_A = T_B$. Similarly, if B and C remain in equilibrium when brought into thermal contact, then $T_B = T_C$. But, we can then conclude that $T_A = T_C$, so systems A and C will also remain in equilibrium when brought into thermal contact. Thus, we arrive at the following statement, which is sometimes called the *zeroth law of thermodynamics*:

If two systems are separately in thermal equilibrium with a third system then they must also be in thermal equilibrium with one another.

Let us test our scheme out on a monatomic ideal gas. We saw in Section 5.4.4 that the number of accessible states of an ideal monatomic gas consisting of N particles is

$$\Omega(U, V) = B V^N U^{3N/2}, \quad (5.318)$$

where U is the internal energy, V the volume, and B is a constant that is independent of U and V . According to the previous two equations, the thermodynamic temperature of such a gas is

$$\frac{1}{k_B T} = \frac{3N}{2U}. \quad (5.319)$$

However, $N = \nu N_A$, where ν is the number of moles of molecules in the gas, and N_A is Avogadro’s number. The previous equation can be rearranged to give

$$U = \frac{3}{2} \nu R T, \quad (5.320)$$

because $R = k_B N_A$. However, this is the correct expression for the internal energy of a monatomic ideal gas. (See Section 5.2.3.) Hence, it is clear that the thermodynamic temperature defined in Equation (5.317) corresponds to the more familiar absolute temperature associated with an ideal gas.

5.4.7 Boltzmann Probability Distribution

We have gained some understanding of the macroscopic properties of the air in a classroom (say). For instance, we know something about its internal energy and specific heat capacity. How can we

obtain information about the statistical properties of the molecules that make up this air? Consider a specific molecule. It constantly collides with its immediate neighbor molecules, and occasionally bounces off the walls of the room. These interactions “inform” it about the macroscopic state of the air, such as its temperature, pressure, and volume. The statistical distribution of the molecule over its own particular internal states must be consistent with this macroscopic state. In other words, if we have a large group of such molecules with similar statistical distributions then they must be equivalent to air with the appropriate macroscopic properties. So, it ought to be possible to calculate the probability distribution of the molecule over its internal states from a knowledge of these macroscopic properties.

We can think of the interaction of a molecule with the air in a classroom as analogous to the interaction of a small system, A , in thermal contact with a heat reservoir, A' . The air acts like a heat reservoir because its energy fluctuations due to interactions with the molecule are far too small to affect any of its macroscopic parameters. Let us determine the probability, P_r , of finding system A in one particular internal state, r , of energy ϵ_r , when it is thermal equilibrium with the heat reservoir, A' .

As usual, we assume fairly weak interaction between A and A' , so that the energies of these two systems are additive. The energy of A is not known at this stage. In fact, only the total internal energy of the combined system, $A^{(0)} = A + A'$, is known. Suppose that the total internal energy lies in the range $U^{(0)}$ to $U^{(0)} + \delta U$. The overall internal energy is constant in time, because $A^{(0)}$ is assumed to be an isolated system, so

$$\epsilon_r + U' = U^{(0)}, \quad (5.321)$$

where U' denotes the internal energy of the reservoir A' . Let $\Omega'(U')$ be the number of accessible states of the reservoir when its internal energy lies in the range U' to $U' + \delta U$. Clearly, if system A has an energy ϵ_r , then the reservoir A' must have an energy close to $U' = U^{(0)} - \epsilon_r$. Hence, because A is in one definite state (i.e., state r), and the total number of states accessible to A' is $\Omega'(U^{(0)} - \epsilon_r)$, it follows that the total number of states accessible to the combined system is simply $\Omega'(U^{(0)} - \epsilon_r)$. The principle of equal a priori probabilities tells us the probability of occurrence of a particular situation is proportional to the number of accessible states. Thus,

$$P_r = C' \Omega'(U^{(0)} - \epsilon_r), \quad (5.322)$$

where C' is a constant of proportionality that is independent of r . This constant can be determined by the normalization condition

$$\sum_r P_r = 1, \quad (5.323)$$

where the sum is over all possible states of system A , irrespective of their energy. [See Equation (5.3).]

Let us now make use of the fact that system A is far smaller than system A' . It follows that $\epsilon_r \ll U^{(0)}$, so the slowly-varying logarithm of P_r can be Taylor expanded about $U' = U^{(0)}$. Thus,

$$\ln P_r = \ln C' + \ln \Omega'(U^{(0)}) - \left[\frac{\partial \ln \Omega'}{\partial U'} \right]_0 \epsilon_r + \dots \quad (5.324)$$

Note that we must expand $\ln P_r$, rather than P_r itself, because the latter function varies so rapidly with energy that the radius of convergence of its Taylor series is too small for the series to be of any practical use. The higher-order terms in Equation (5.324) can be safely neglected, because $\epsilon_r \ll U^{(0)}$. Now, the derivative

$$\left[\frac{\partial \ln \Omega'}{\partial U'} \right]_0 \equiv \beta \quad (5.325)$$

is evaluated at the fixed energy $U' = U^{(0)}$, and is, thus, a constant, independent of the energy, ϵ_r , of A . In fact, we know, from the previous section, that this derivative is just the temperature parameter $\beta = (k_B T)^{-1}$ characterizing the heat reservoir A' . Here, T is the absolute temperature of the reservoir. Hence, Equation (5.324) becomes

$$\ln P_r = \ln C' + \ln \Omega'(U^{(0)}) - \frac{\epsilon_r}{k_B T}, \quad (5.326)$$

giving

$$P_r = C \exp\left(-\frac{\epsilon_r}{k_B T}\right), \quad (5.327)$$

where C is a constant independent of r . The parameter C is determined by the normalization condition, (5.323), which gives

$$C^{-1} = \sum_r \exp\left(-\frac{\epsilon_r}{k_B T}\right). \quad (5.328)$$

We conclude that the probability of a measurement of the energy of some system A , that is in thermal equilibrium with a heat reservoir of temperature T , yielding the result ϵ_r is

$$P_r = \frac{\exp(-\epsilon_r/k_B T)}{\sum_r \exp(-\epsilon_r/k_B T)}. \quad (5.329)$$

This probability distribution is known as the *Boltzmann probability distribution*.

5.5 Applications of Statistical Mechanics

5.5.1 Two-State System

Consider a microscopic system (such as an atom) that possesses two quantum states, labelled 1 and 2. Let the lower energy state, 1 (i.e., the ground state), have energy 0, and let the higher energy state (i.e., the excited state), 2, have energy Δ , where $\Delta > 0$.

Suppose that the microscopic system is in thermal equilibrium with a heat reservoir of temperature T . According to the Boltzmann distribution, (5.329), the probability the system is found in state i is

$$P_i = \frac{\exp(\epsilon_i/k_B T)}{\exp(\epsilon_1/k_B T) + \exp(-\epsilon_2/k_B T)}, \quad (5.330)$$

where $i = 1, 2$. In particular, given that $\epsilon_1 = 0$ and $\epsilon_2 = \Delta$, we find that

$$P_1 = \frac{1}{1 + \exp(-\Delta/k_B T)}, \quad (5.331)$$

$$P_2 = \frac{1}{1 + \exp(\Delta/k_B T)}. \quad (5.332)$$

Note that $P_1 + P_2 = 1$. Thus, at low temperatures, $k_B T \ll \Delta$, we obtain $P_1 \rightarrow 1$ and $P_2 \rightarrow 0$. In other words, at low temperatures, the system is certain to be found in its ground state, and has no chance of being found in its excited state. On the other hand, at high temperatures, $k_B T \gg \Delta$, we obtain $P_1 = P_2 = 1/2$. In other words, at high temperatures, the microscopic system is equally likely to be found in its ground state or in its excited state. Finally, the mean energy of the microscopic system is (see Section 5.1.3)

$$\langle E \rangle = P_1 \epsilon_1 + P_2 \epsilon_2 = \frac{\Delta}{1 + \exp(\Delta/k_B T)}. \quad (5.333)$$

Note that there is no temperature at which it is possible to get a *population inversion*; that is, $P_2 > P_1$. In fact, lasers, which require a population inversion in order to operate, are not in thermal equilibrium.

Suppose that we have a macroscopic system consisting of N identical two-state microscopic systems of the type that we have just discussed. The internal energy of the macroscopic system is

$$U = N \langle \epsilon \rangle = \frac{N \Delta}{1 + \exp(\Delta/k_B T)}. \quad (5.334)$$

Moreover, the specific heat capacity of the macroscopic system at constant volume is (see Section 5.2.3)

$$C_V = \left(\frac{\partial U}{\partial T} \right)_{V,N} = \frac{N \Delta^2}{k_B T^2} \frac{\exp(\Delta/k_B T)}{[1 + \exp(\Delta/k_B T)]^2}. \quad (5.335)$$

The previous two equations yield

$$\frac{U}{N k_B T_c} = \frac{\exp(-T_c/T)}{\cosh(T_c/T)}, \quad (5.336)$$

$$\frac{C_V}{N k_B} = \frac{(T_c/T)^2}{\cosh^2(T_c/T)}, \quad (5.337)$$

where $T_c = \Delta/(2k_B)$. Figure 5.2 illustrates how U and C_V vary with temperature. The peak in the heat capacity is known as the *Schottky anomaly*, and is associated with the absorption of energy from the heat reservoir as the temperature exceeds the critical temperature required for the constituent microscopic systems to be excited from their ground states.

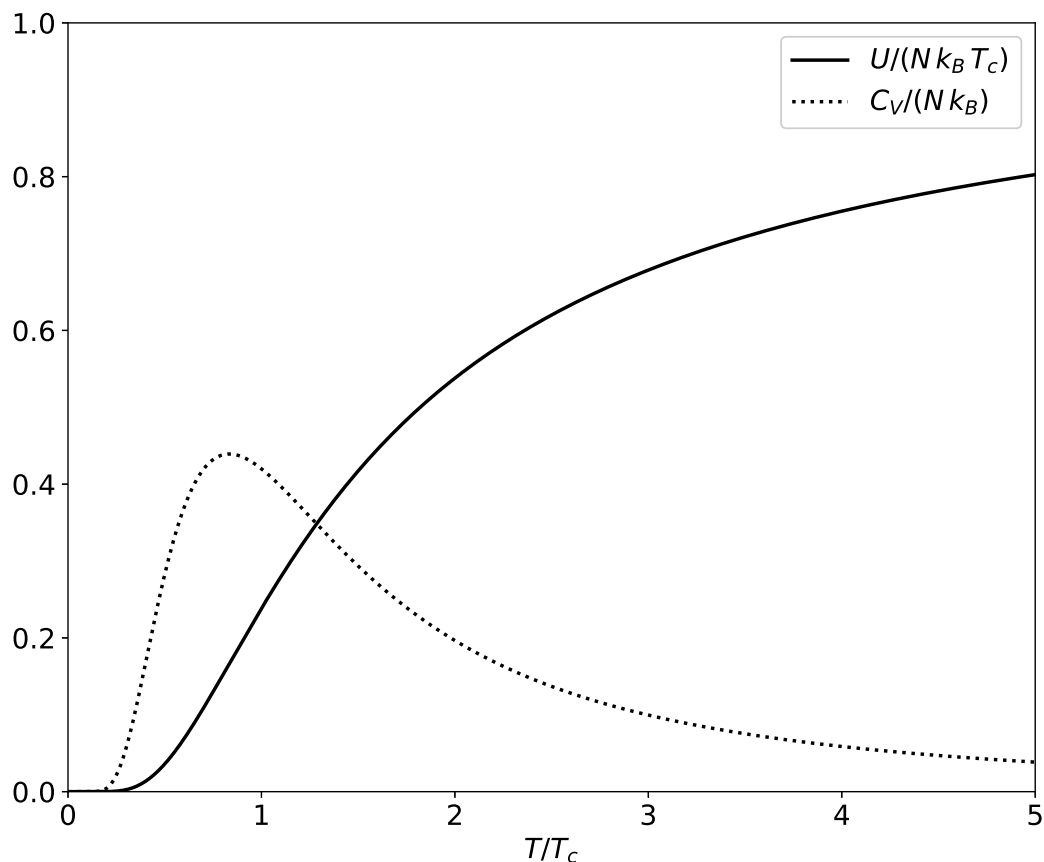


Figure 5.2: Internal energy and specific heat capacity of a two-state system as a function of the temperature.

5.5.2 Spin-1/2 Paramagnetism

As a specific example of a two-state system, consider a substance whose constituent atoms contain only one unpaired electron (with zero orbital angular momentum). Such atoms have spin $1/2$ [i.e., their spin angular momentum is $(1/2)\hbar$], and consequently possess an intrinsic magnetic moment, μ . According to quantum mechanics, the magnetic moment of a spin- $1/2$ atom can point either parallel or antiparallel to an external magnetic field, \mathbf{B} . Let us determine the mean magnetic moment (parallel to \mathbf{B}), $\langle\mu_{\parallel}\rangle$, of the constituent atoms of the substance when its absolute temperature is T . We shall assume, for the sake of simplicity, that each atom only interacts weakly with its neighboring atoms. This enables us to focus attention on a single atom, and to treat the remaining atoms as a heat reservoir at temperature T .

Our atom can be in one of two possible states. Namely, the (+) state in which its spin points up (i.e., parallel to \mathbf{B}), or the (−) state in which its spin points down (i.e., antiparallel to \mathbf{B}). In the (+)

state, the atomic magnetic moment is parallel to the magnetic field, so that $\mu_{\parallel} = \mu$. The magnetic energy of the atom is $\epsilon_+ = -\mu B$. In the (-) state, the atomic magnetic moment is antiparallel to the magnetic field, so that $\mu_{\parallel} = -\mu$. The magnetic energy of the atom is $\epsilon_- = \mu B$.

According to the Boltzmann distribution, (5.329), the probability of finding the atom in the (+) state is

$$P_+ = \frac{\exp(-\epsilon_+/k_B T)}{\exp(-\epsilon_+/k_B T) + \exp(-\epsilon_-/k_B T)} = \frac{\exp(\mu B/k_B T)}{\exp(\mu B/k_B T) + \exp(-\mu B/k_B T)}, \quad (5.338)$$

Likewise, the probability of finding the atom in the (-) state is

$$P_- = \frac{\exp(-\epsilon_-/k_B T)}{\exp(-\epsilon_+/k_B T) + \exp(-\epsilon_-/k_B T)} = \frac{\exp(-\mu B/k_B T)}{\exp(\mu B/k_B T) + \exp(-\mu B/k_B T)}. \quad (5.339)$$

Clearly, the most probable state is the state with the lower energy [i.e., the (+) state]. Thus, the mean magnetic moment points in the direction of the magnetic field (i.e., the atomic spin is more likely to point parallel to the field than antiparallel).

It is apparent that the critical parameter in a paramagnetic system is

$$y = \frac{\mu B}{k_B T}. \quad (5.340)$$

This dimensionless parameter measures the ratio of the typical magnetic energy of the atom, μB , to its typical thermal energy, $k_B T$. If the thermal energy greatly exceeds the magnetic energy then $y \ll 1$, and the probability that the atomic moment points parallel to the magnetic field is about the same as the probability that it points antiparallel. In this situation, we expect the mean atomic moment to be small, so that $\langle \mu_{\parallel} \rangle \simeq 0$. On the other hand, if the magnetic energy greatly exceeds the thermal energy then $y \gg 1$, and the atomic moment is far more likely to be directed parallel to the magnetic field than antiparallel. In this situation, we expect $\langle \mu_{\parallel} \rangle \simeq \mu$.

Let us calculate the mean atomic moment, $\langle \mu_{\parallel} \rangle$. The usual definition of a mean value gives (see Section 5.1.3)

$$\langle \mu_{\parallel} \rangle = \frac{P_+ \mu + P_- (-\mu)}{P_+ + P_-} = \mu \left[\frac{\exp(\mu B/k_B T) - \exp(-\mu B/k_B T)}{\exp(\mu B/k_B T) + \exp(-\mu B/k_B T)} \right]. \quad (5.341)$$

This can also be written

$$\langle \mu_{\parallel} \rangle = \mu \tanh \left(\frac{\mu B}{k_B T} \right). \quad (5.342)$$

For small arguments, $y \ll 1$,

$$\tanh y \simeq y - \frac{y^3}{3} + \dots, \quad (5.343)$$

whereas for large arguments, $y \gg 1$,

$$\tanh y \simeq 1. \quad (5.344)$$

It follows that at comparatively high temperatures, $k_B T \gg \mu B$,

$$\langle \mu_{\parallel} \rangle \simeq \frac{\mu^2 B}{k_B T}, \quad (5.345)$$

whereas at comparatively low temperatures, $k_B T \ll \mu B$,

$$\langle \mu_{\parallel} \rangle \simeq \mu. \quad (5.346)$$

Suppose that the substance contains N_0 atoms per unit volume. The *magnetization* is defined as the mean magnetic moment per unit volume, and is given by

$$\langle M_{\parallel} \rangle = N_0 \langle \mu_{\parallel} \rangle. \quad (5.347)$$

At high temperatures, $k_B T \gg \mu B$, the mean magnetic moment, and, hence, the magnetization, is proportional to the applied magnetic field, so we can write

$$\langle M_{\parallel} \rangle \simeq \chi \frac{B}{\mu_0}, \quad (5.348)$$

where χ is a dimensionless constant of proportionality known as the *magnetic susceptibility*, and μ_0 the magnetic permeability of free space. It is clear that the magnetic susceptibility of a spin-1/2 paramagnetic substance takes the form

$$\chi = \frac{N_0 \mu_0 \mu^2}{k_B T}. \quad (5.349)$$

The fact that $\chi \propto T^{-1}$ is known as *Curie's law*, because it was discovered experimentally by Pierre Curie at the end of the nineteenth century. At low temperatures, $k_B T \ll \mu B$,

$$\langle M_{\parallel} \rangle \rightarrow N_0 \mu, \quad (5.350)$$

so the magnetization becomes independent of the applied field. This corresponds to the maximum possible magnetization, in which all atomic moments are aligned parallel to the field. The breakdown of the $\langle M_{\parallel} \rangle \propto B$ law at low temperatures (or high magnetic fields) is known as *saturation*.

5.5.3 Adiabatic Demagnetization

Suppose that we take the spin-1/2 paramagnetic system discussed in the previous section, and thermally isolate it from its surroundings. In this case, the numbers of atoms in the spin-up and spin-down states cannot change, because the system is unable to get rid of excess energy. In other words, the ratio of the number of atoms in the spin-up state to the number of atoms in the spin-down state,

$$\frac{N_+}{N_-} = \frac{P_+}{P_-} = \exp\left(\frac{2\mu B}{k_B T}\right), \quad (5.351)$$

is fixed. Under these so-called *adiabatic* conditions, we find that $T \propto B$. This is known as the *magnetocaloric effect*.

The magnetocaloric effect is the basis of a method of cooling atomic systems down to very low temperatures that is known as *adiabatic demagnetization*. In this scheme, the sample is initially in thermal contact with liquid helium at 0.8 K. The sample is then magnetized. In the process, heat is given off by the sample, and is conducted away by the liquid helium. Next, the sample is thermally isolated by pumping out the liquid helium. Finally, the sample is demagnetized, leading to a reduction in its temperature via the magnetocaloric effect. Temperatures as low as 10^{-8} K have been achieved by this method.

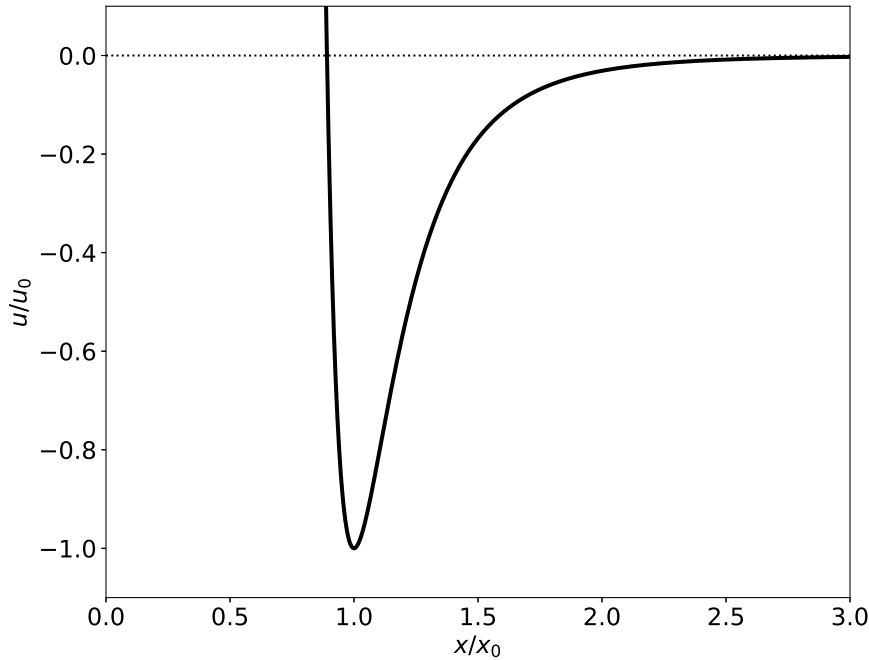


Figure 5.3: The Lennard-Jones potential.

5.5.4 Thermal Expansion

The interaction of electrically neutral atoms can be modeled using the *Lennard-Jones potential*,

$$u(x) = u_0 \left[\left(\frac{x_0}{x} \right)^{12} - 2 \left(\frac{x_0}{x} \right)^6 \right], \quad (5.352)$$

where $u(x)$ is the potential energy of a pair of atoms when they are a distance x apart, and u_0 and x_0 are positive constants. See Figure 5.3. We can think of a given atom in a solid made up of neutral atoms as a microscopic system interacting with a heat reservoir that consists of all of the other atoms. Let T be the temperature of the reservoir. Let us treat the problem classically, which is equivalent to assuming that the temperature is sufficiently high that an atom moving in the previous potential is distributed over a large number of different quantum states. According to a straightforward generalization of the Boltzmann distribution, (5.329), the mean value of x is

$$\langle x \rangle = \frac{\int_0^\infty \exp[-u(x)/(k_B T)] x dx}{\int_0^\infty \exp[-u(x)/(k_B T)] dx}. \quad (5.353)$$

Let us assume that the temperature is sufficiently low that an atom is only likely to be found relatively close to the bottom of the potential well, $x = x_0$. We can Taylor expand the potential about $x = x_0$ to give

$$u(x) = u_0 + u'_0 (x - x_0) + \frac{1}{2} u''_0 (x - x_0)^2 + \frac{1}{6} u'''_0 (x - x_0)^3 + \cdots, \quad (5.354)$$

where

$$u_0 = u(x_0), \quad (5.355)$$

$$u'_0 = \frac{du(x_0)}{dx} = 0, \quad (5.356)$$

$$u''_0 = \frac{d^2u(x_0)}{dx^2} = \frac{72 u_0}{x_0^2}, \quad (5.357)$$

$$u'''_0 = \frac{d^3u(x_0)}{dx^3} = -\frac{1512 u_0}{x_0^3}. \quad (5.358)$$

Thus, Equation (5.353) gives

$$\langle x \rangle = \frac{\int_{-\infty}^{\infty} \exp[-u'_0 y^2 / (2 k_B T)] \exp[-u'''_0 y^3 / (6 k_B T)] (x_0 + y) dy}{\int_{-\infty}^{\infty} \exp[-u'_0 y^2 / (2 k_B T)] \exp[-u'''_0 y^3 / (6 k_B T)] dy}, \quad (5.359)$$

where $y = x - x_0$, and we can safely replace the lower limits of integration by $-\infty$, in the integrals on the right-hand side of the previous expression, because we are assuming that the atom is very unlikely to be found a large distance from the bottom of the potential. Let us further assume that $|u'''_0 y^3 / (6 k_B T)| \ll 1$. In this case, we can write

$$\langle x \rangle \simeq x_0 + \frac{\int_{-\infty}^{\infty} \exp[-u'_0 y^2 / (2 k_B T)] [1 - u'''_0 y^3 / (6 k_B T)] y dy}{\int_{-\infty}^{\infty} \exp[-u'_0 y^2 / (2 k_B T)] [1 - u'''_0 y^3 / (6 k_B T)] dy}. \quad (5.360)$$

Now, $\exp[-u'_0 y^2 / (2 k_B T)]$ and y^4 are even functions of y , whereas y and y^3 are odd functions. In general, an integral over all y of the product of an even and an odd function averages to zero, whereas an integral of the product of two even functions does not. Hence, the previous equation simplifies to give

$$\begin{aligned} \langle x \rangle &= x_0 - \frac{u'''_0}{6 k_B T} \frac{\int_{-\infty}^{\infty} y^4 \exp[-u'_0 y^2 / (2 k_B T)] dy}{\int_{-\infty}^{\infty} \exp[-u'_0 y^2 / (2 k_B T)] dy} \\ &= x_0 + \left(\frac{-u'''_0}{6 k_B T} \right) \left(\frac{2 k_B T}{u'_0} \right)^2 \frac{\int_{-\infty}^{\infty} z^4 \exp(-z^2) dz}{\int_{-\infty}^{\infty} \exp(-z^2) dz}. \end{aligned} \quad (5.361)$$

However, $\int_{-\infty}^{\infty} z^4 \exp(-z^2) dz = (3/4) \pi^{1/2}$, and $\int_{-\infty}^{\infty} \exp(-z^2) dz = \pi^{1/2}$, so we get

$$\langle x \rangle = x_0 + \frac{(-u'''_0) k_B T}{2 (u'_0)^2}. \quad (5.362)$$

Note that $\langle x \rangle$ increases linearly with increasing temperature.

The *coefficient of linear thermal expansion* of a solid is defined

$$\alpha = \frac{1}{\langle x \rangle} \frac{d\langle x \rangle}{dT}, \quad (5.363)$$

where $\langle x \rangle$ is the mean distance between nearest neighbor atoms. The previous two equations yield

$$\alpha = \frac{(-u_0'') k_B}{2 x_0 (u_0')^2} = \frac{7 k_B}{48 u_0}, \quad (5.364)$$

where use has been made of Equations (5.357) and (5.358).

For solid argon at 80 K, $x_0 = 3.9 \times 10^{-10}$ m, and $u_0 = 0.010$ eV. Hence, we deduce that

$$\alpha = 1.3 \times 10^{-3} \text{ K}^{-1}. \quad (5.365)$$

The measured value of α is about $2 \times 10^{-3} \text{ K}^{-1}$.

5.5.5 Equipartition Theorem

The internal energy of a monatomic ideal gas containing N atoms is $(3/2) N k_B T$. (See Section 5.2.3.) This implies that each atom possess, on average, $(3/2) k_B T$ units of energy. Monatomic particles have only three translational degrees of freedom, corresponding to their motion in three dimensions. They possess no internal rotational or vibrational degrees of freedom. Thus, the mean energy per degree of freedom in a monatomic ideal gas is $(1/2) k_B T$. In fact, this is a special case of a more general result. Let us now try to prove this result.

Suppose that the energy of a system is determined by f coordinates, q_k , and f corresponding momenta, p_k , so that

$$E = E(q_1, \dots, q_f, p_1, \dots, p_f). \quad (5.366)$$

Suppose further that:

1. The total energy splits additively into the form

$$E = \epsilon_i(p_i) + E'(q_1, \dots, p_f), \quad (5.367)$$

where ϵ_i involves only one variable, p_i , and the remaining part, E' , does not depend on p_i .

2. The function ϵ_i is quadratic in p_i , so that

$$\epsilon_i(p_i) = b p_i^2, \quad (5.368)$$

where b is a constant.

The most common situation in which the previous assumptions are valid is where p_i is a momentum. This is because the kinetic energy is usually a quadratic function of each momentum component, whereas the potential energy does not involve the momenta at all. However, if a coordinate, q_i , were to satisfy assumptions 1 and 2 then the theorem that we are about to establish would hold just as well.

What is the mean value of ϵ_i in thermal equilibrium if conditions 1 and 2 are satisfied? If the system is in equilibrium at absolute temperature T then it is distributed according to the Boltzmann

probability distribution. (See Section 5.4.7.) In the classical approximation, the mean value of ϵ_i is expressed in terms of integrals over all phase-space (see Section 5.4.4):

$$\langle \epsilon_i \rangle = \frac{\int_{-\infty}^{\infty} \exp[-E(q_1, \dots, p_f)/k_B T] \epsilon_i dq_1 \cdots dp_f}{\int_{-\infty}^{\infty} \exp[-E(q_1, \dots, p_f)/k_B T] dq_1 \cdots dp_f}. \quad (5.369)$$

Condition 1 gives

$$\begin{aligned} \langle \epsilon_i \rangle &= \frac{\int_{-\infty}^{\infty} \exp[-(\epsilon_i + E')/k_B T] \epsilon_i dq_1 \cdots dp_f}{\int_{-\infty}^{\infty} \exp[-(\epsilon_i + E')/k_B T] dq_1 \cdots dp_f} \\ &= \frac{\int_{-\infty}^{\infty} \exp(-\epsilon_i/k_B T) \epsilon_i dp_i \int_{-\infty}^{\infty} \exp(-E'/k_B T) dq_1 \cdots dp_f}{\int_{-\infty}^{\infty} \exp(-\epsilon_i/k_B T) dp_i \int_{-\infty}^{\infty} \exp(-E'/k_B T) dq_1 \cdots dp_f}, \end{aligned} \quad (5.370)$$

where use has been made of the multiplicative property of the exponential function, and where the final integrals in both the numerator and denominator extend over all variables, q_k and p_k , except for p_i . These integrals are equal and, thus, cancel. Hence,

$$\langle \epsilon_i \rangle = \frac{\int_{-\infty}^{\infty} \exp(-\epsilon_i/k_B T) \epsilon_i dp_i}{\int_{-\infty}^{\infty} \exp(-\epsilon_i/k_B T) dp_i}. \quad (5.371)$$

This expression can be simplified further because, writing $\beta = 1/k_B T$,

$$\int_{-\infty}^{\infty} \exp(-\beta \epsilon_i) \epsilon_i dp_i \equiv -\frac{\partial}{\partial \beta} \left[\int_{-\infty}^{\infty} \exp(-\beta \epsilon_i) dp_i \right], \quad (5.372)$$

so

$$\langle \epsilon_i \rangle = -\frac{\partial}{\partial \beta} \ln \left[\int_{-\infty}^{\infty} \exp(-\beta \epsilon_i) dp_i \right]. \quad (5.373)$$

According to condition 2,

$$\int_{-\infty}^{\infty} \exp(-\beta \epsilon_i) dp_i = \int_{-\infty}^{\infty} \exp(-\beta b p_i^2) dp_i = \frac{1}{\sqrt{\beta}} \int_{-\infty}^{\infty} \exp(-b y^2) dy, \quad (5.374)$$

where $y = \sqrt{\beta} p_i$. Thus,

$$\ln \int_{-\infty}^{\infty} \exp(-\beta \epsilon_i) dp_i = -\frac{1}{2} \ln \beta + \ln \int_{-\infty}^{\infty} \exp(-b y^2) dy. \quad (5.375)$$

Note that the integral on the right-hand side is independent of β . It follows from Equation (5.373) that

$$\langle \epsilon_i \rangle = -\frac{\partial}{\partial \beta} \left(-\frac{1}{2} \ln \beta \right) = \frac{1}{2\beta}, \quad (5.376)$$

giving

$$\langle \epsilon_i \rangle = \frac{1}{2} k_B T. \quad (5.377)$$

This result is known as the *equipartition theorem*. It states that the mean value of every independent quadratic term in the energy is equal to $(1/2) k_B T$. If all terms in the energy are quadratic then the mean energy is spread equally over all degrees of freedom. (Hence, the name “equipartition.”)

5.5.6 Harmonic Oscillator

Our proof of the equipartition theorem depends crucially on the classical approximation. To see how quantum effects modify this result, let us examine a particularly simple system that we know how to analyze using both classical and quantum physics; namely, a simple harmonic oscillator. Consider a one-dimensional harmonic oscillator in equilibrium with a heat reservoir held at absolute temperature T . The energy of the oscillator is given by

$$\epsilon = \frac{p^2}{2m} + \frac{1}{2} \kappa x^2, \quad (5.378)$$

where the first term on the right-hand side is the kinetic energy, involving the momentum, p , and the mass, m , and the second term is the potential energy, involving the displacement, x , and the force constant, κ . Each of these terms is quadratic in the respective variable. So, in the classical approximation, the equipartition theorem yields:

$$\frac{\langle p^2 \rangle}{2m} = \frac{1}{2} k_B T, \quad (5.379)$$

$$\frac{1}{2} \kappa \langle x^2 \rangle = \frac{1}{2} k_B T. \quad (5.380)$$

That is, the mean kinetic energy of the oscillator is equal to the mean potential energy, which equals $(1/2) k_B T$. It follows that the mean total energy is

$$\langle \epsilon \rangle = \frac{1}{2} k_B T + \frac{1}{2} k_B T = k_B T. \quad (5.381)$$

According to quantum mechanics (see Section 4.3.7), the energy levels of a harmonic oscillator are equally spaced, and satisfy

$$\epsilon_n = \left(\frac{1}{2} + n \right) \hbar \omega, \quad (5.382)$$

where n is a non-negative integer, and

$$\omega = \sqrt{\frac{\kappa}{m}}. \quad (5.383)$$

Making use of the Boltzmann distribution, (5.329), the mean value of the quantum number n for such an oscillator is

$$\langle n \rangle = \frac{\sum_{n=0, \infty} \exp[-(n + 1/2) \hbar \omega / k_B T] n}{\sum_{n=0, \infty} \exp[-(n + 1/2) \hbar \omega / k_B T]} = \frac{\sum_{n=0, \infty} n x^n}{\sum_{n=0, \infty} x^n} = x \frac{d}{dx} \left[\ln \sum_{n=0, \infty} x^n \right], \quad (5.384)$$

where $x = \exp(-\hbar \omega / k_B T)$. Now,

$$\sum_{n=0, \infty} x^n = \frac{1}{1 - x}, \quad (5.385)$$

so

$$\ln \sum_{n=0, \infty} x^n = -\ln(1 - x), \quad (5.386)$$

and

$$x \frac{d}{dx} \left[\ln \sum_{n=0, \infty} x^n \right] = \frac{x}{1-x}. \quad (5.387)$$

Hence, we deduce that

$$\langle n \rangle = \frac{1}{\exp(\hbar \omega / k_B T) - 1}. \quad (5.388)$$

Thus, the mean energy of the oscillator,

$$\langle \epsilon \rangle = \left(\frac{1}{2} + \langle n \rangle \right) \hbar \omega, \quad (5.389)$$

takes the form

$$\langle \epsilon \rangle = \left[\frac{1}{2} + \frac{1}{\exp(\hbar \omega / k_B T) - 1} \right] \hbar \omega. \quad (5.390)$$

Consider the limit

$$\frac{\hbar \omega}{k_B T} \ll 1, \quad (5.391)$$

in which the thermal energy, $k_B T$, is large compared to the separation, $\hbar \omega$, between successive energy levels. In this limit,

$$\exp\left(\frac{\hbar \omega}{k_B T}\right) \simeq 1 + \frac{\hbar \omega}{k_B T}, \quad (5.392)$$

so

$$\langle \epsilon \rangle \simeq \left(\frac{1}{2} + \frac{k_B T}{\hbar \omega} \right) \hbar \omega \simeq \left(\frac{k_B T}{\hbar \omega} \right) \hbar \omega, \quad (5.393)$$

giving

$$\langle \epsilon \rangle \simeq k_B T. \quad (5.394)$$

Thus, the classical result, (5.381), holds whenever the thermal energy greatly exceeds the typical spacing between quantum energy levels.

Consider the limit

$$\frac{\hbar \omega}{k_B T} \gg 1, \quad (5.395)$$

in which the thermal energy is small compared to the separation between the energy levels. In this limit,

$$\exp\left(\frac{\hbar \omega}{k_B T}\right) \gg 1, \quad (5.396)$$

and so

$$\langle \epsilon \rangle \simeq \left[\frac{1}{2} + \exp\left(-\frac{\hbar \omega}{k_B T}\right) \right] \hbar \omega \simeq \frac{1}{2} \hbar \omega. \quad (5.397)$$

Thus, if the thermal energy is much less than the spacing between quantum states then the mean energy approaches that of the ground state. Clearly, the equipartition theorem is only valid in the former limit, where $k_B T \gg \hbar \omega$, and the oscillator possess sufficient thermal energy to explore many of its possible quantum states.

5.5.7 Specific Heat Capacities

Classical physics, in the guise of the equipartition theorem, tells us that each independent degree of freedom associated with a quadratic term in the energy possesses an average energy $(1/2)k_B T$ in thermal equilibrium at temperature T . Consider a substance made up of N molecules. Every molecular degree of freedom contributes $(1/2)Nk_B T$, or $(1/2)\nu RT$, to the mean internal energy of the substance (with the tacit proviso that each degree of freedom is associated with a quadratic term in the energy). Thus, the contribution to the molar heat capacity at constant volume is

$$\frac{1}{\nu} \left(\frac{\partial U}{\partial T} \right)_V = \frac{1}{\nu} \frac{\partial[(1/2)\nu RT]}{\partial T} = \frac{1}{2} R, \quad (5.398)$$

per molecular degree of freedom. The total classical heat capacity is therefore

$$c_V = \frac{g}{2} R, \quad (5.399)$$

where g is the number of molecular degrees of freedom.

As we have seen, the equipartition theorem (and the whole classical approximation) is only valid when the typical thermal energy, $k_B T$, greatly exceeds the spacing between quantum energy levels. Suppose that the temperature is sufficiently low that this condition is not satisfied for one particular molecular degree of freedom. In fact, suppose that $k_B T$ is much less than the spacing between the energy levels. In this situation, the degree of freedom only contributes the ground-state energy, E_0 (say) to the mean energy of the molecule. Now, the ground-state energy can be a quite complicated function of the internal properties of the molecule, but is certainly not a function of the temperature, because this is a collective property of all molecules. It follows that the contribution to the molar heat capacity is zero. Thus, if $k_B T$ is much less than the spacing between the energy levels then the degree of freedom contributes nothing at all to the molar heat capacity. We say that this particular degree of freedom is “frozen out.” Clearly, at very low temperatures, just about all degrees of freedom are frozen out. As the temperature is gradually increased, degrees of freedom successively kick in, and eventually contribute their full $(1/2)R$ to the molar heat capacity, as $k_B T$ approaches, and then greatly exceeds, the spacing between their quantum energy levels. We can use these simple ideas to explain the behaviors of most experimental heat capacities.

To make further progress, we need to estimate the typical spacing between the quantum energy levels associated with various degrees of freedom. We can do this by observing the frequency of the electromagnetic radiation emitted and absorbed during transitions between these energy levels. If the typical spacing between energy levels is ΔE then transitions between the various levels are associated with photons of frequency ν , where $h\nu = \Delta E$. (Here, h is Planck’s constant.) We can define an *effective temperature* of the radiation via $h\nu = k_B T_{\text{rad}}$. If $T \gg T_{\text{rad}}$ then $k_B T \gg \Delta E$, and the degree of freedom makes its full contribution to the heat capacity. On the other hand, if $T \ll T_{\text{rad}}$ then $k_B T \ll \Delta E$, and the degree of freedom is frozen out. Table 5.1 lists the “temperatures” of various different types of radiation. It is clear that degrees of freedom that give rise to emission or absorption of radio or microwave radiation contribute their full $(1/2)R$ to the molar heat capacity at room temperature. On the other hand, degrees of freedom that give rise to emission or absorption in the visible, ultraviolet, X-ray, or γ -ray regions of the electromagnetic spectrum are frozen out

Radiation type	Frequency (hz)	$T_{\text{rad}}(\text{K})$
Radio	$< 10^9$	< 0.05
Microwave	$10^9 - 10^{11}$	$0.05 - 5$
Infrared	$10^{11} - 10^{14}$	$5 - 5000$
Visible	5×10^{14}	2×10^4
Ultraviolet	$10^{15} - 10^{17}$	$5 \times 10^4 - 5 \times 10^6$
X-ray	$10^{17} - 10^{20}$	$5 \times 10^6 - 5 \times 10^9$
γ -ray	$> 10^{20}$	$> 5 \times 10^9$

Table 5.1: Effective “temperatures” of various different types of electromagnetic radiation.

at room temperature. Degrees of freedom that emit or absorb infrared radiation are on the border line.

5.5.8 Specific Heats of Gases

Let us now investigate the specific heats of gases. Consider, first of all, translational degrees of freedom. Every molecule in a gas is free to move in three dimensions. If one particular molecule has mass m and momentum $\mathbf{p} = m \mathbf{v}$ then its kinetic energy of translation is

$$K = \frac{1}{2m} (p_x^2 + p_y^2 + p_z^2). \quad (5.400)$$

The kinetic energy of other molecules does not involve the momentum, \mathbf{p} , of this particular molecule. Moreover, the potential energy of interaction between molecules depends only on their position coordinates, and is, thus, independent of \mathbf{p} . Any internal rotational, vibrational, electronic, or nuclear degrees of freedom of the molecule also do not involve \mathbf{p} . Hence, the essential conditions of the equipartition theorem are satisfied. (At least, in the classical approximation.) Because Equation (5.400) contains three independent quadratic terms, there are clearly three degrees of freedom associated with translation (one for each dimension of space), so the translational contribution to the molar heat capacity of gases is

$$(c_V)_{\text{translation}} = \frac{3}{2} R. \quad (5.401)$$

Suppose that our gas is contained in a cubic enclosure of dimensions a . According to Schrödinger’s equation, the quantized translational energy levels of an individual molecule are given by

$$E = \frac{\hbar^2 \pi^2}{2m a^2} (n_x^2 + n_y^2 + n_z^2), \quad (5.402)$$

where n_x , n_y , and n_z are positive-integer quantum numbers. (See Section 4.4.2.) Clearly, the spacing between the energy levels can be made arbitrarily small by increasing the size of the enclosure. This implies that translational degrees of freedom can be treated classically, so that

Equation (5.401) is always valid. (Except very close to absolute zero.) We conclude that all gases possess a minimum molar heat capacity of $(3/2)R$ due to the translational degrees of freedom of their constituent molecules.

The electronic degrees of freedom of gas molecules (i.e., the possible configurations of electrons orbiting the atomic nuclei) typically give rise to absorption and emission in the ultraviolet or visible regions of the spectrum. It follows from Table 5.1 that electronic degrees of freedom are frozen out at room temperature. Similarly, nuclear degrees of freedom (i.e., the possible configurations of protons and neutrons in the atomic nuclei) are frozen out because they are associated with absorption and emission in the X-ray and γ -ray regions of the electromagnetic spectrum. In fact, the only additional degrees of freedom that we need worry about for gases are rotational and vibrational degrees of freedom. These typically give rise to absorption lines in the infrared region of the spectrum.

The rotational kinetic energy of a molecule tumbling in space can be written

$$K = \frac{1}{2} \frac{L_x^2}{I_{xx}} + \frac{1}{2} \frac{L_y^2}{I_{yy}} + \frac{1}{2} \frac{L_z^2}{I_{zz}}, \quad (5.403)$$

where the x -, y -, and z -axes are the so called *principal axes of rotation* of the molecule (these are mutually perpendicular), L_x , L_y , and L_z are the angular momenta about these axes, and I_{xx} , I_{yy} , and I_{zz} are the principal moments of inertia about these axes. (See Sections 1.7.2 and 1.7.3.) No other degrees of freedom depend on the angular momenta. Because the kinetic energy of rotation is the sum of three quadratic terms, the rotational contribution to the molar heat capacity of gases is

$$(c_V)_{\text{rotation}} = \frac{3}{2} R, \quad (5.404)$$

according to the equipartition theorem. Note that the typical magnitude of a molecular moment of inertia is md^2 , where m is the molecular mass, and d is the typical interatomic spacing in the molecule. A special case arises if the molecule is linear (e.g., if the molecule is diatomic). In this case, one of the principal axes lies along the line of centers of the atoms. The moment of inertia about this axis is of order $m_e d^2$, where m_e is the electron mass. (See Section 5.3.6.) Because $m_e \ll m$, it follows that the moment of inertia about the line of centers is minuscule compared to the moments of inertia about the other two principal axes. In quantum mechanics, angular momentum is quantized in units of \hbar . The energy levels of a rigid rotator spinning about a principal axis are written

$$E = \frac{\hbar^2}{2I} J(J+1), \quad (5.405)$$

where I is the moment of inertia, and J is a non-negative integer. Note the inverse dependence of the spacing between energy levels on the moment of inertia. It is clear that, for the case of a linear molecule, the rotational degree of freedom associated with spinning along the line of centers of the atoms is frozen out at room temperature, given the very small moment of inertia along this axis, and, hence, the very widely spaced rotational energy levels. Thus, the rotational contribution to the molar heat capacity of a diatomic gas is

$$(c_V)_{\text{rotation}} = R. \quad (5.406)$$

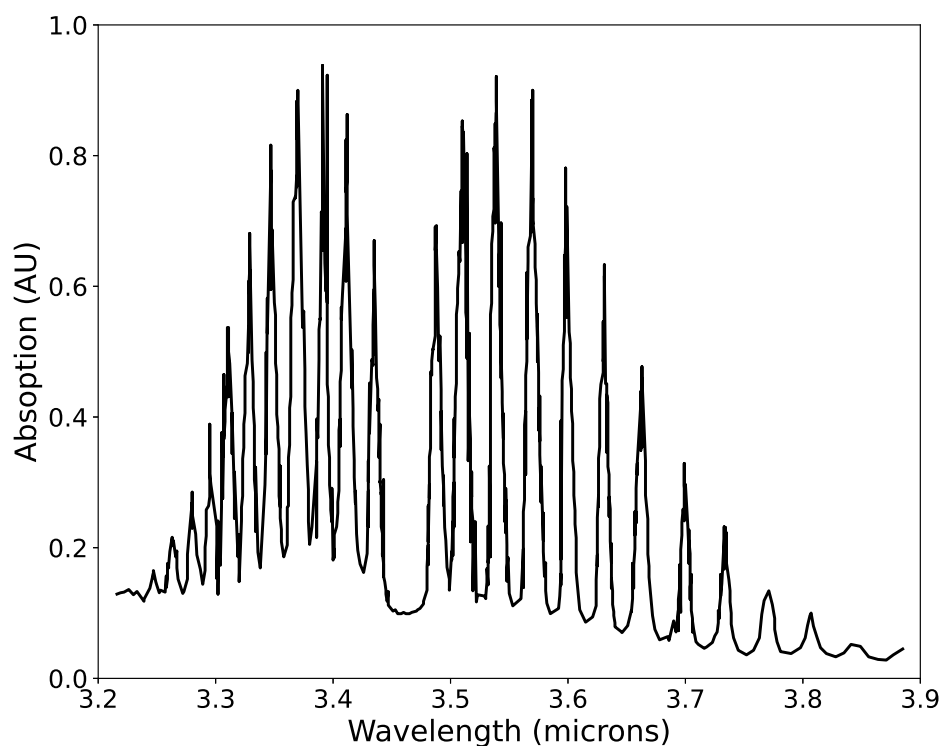


Figure 5.4: The infrared vibration-absorption spectrum of hydrogen chloride gas.

Classically, the vibrational degrees of freedom of a molecule are studied by standard normal mode analysis of the molecular structure. Each normal mode behaves like an independent harmonic oscillator, and, therefore, contributes R to the molar specific heat of the gas [$(1/2)R$ from the kinetic energy of vibration, and $(1/2)R$ from the potential energy of vibration]. A molecule containing n atoms has $n - 1$ normal modes of vibration. For instance, a diatomic molecule has just one normal mode (corresponding to periodic stretching of the bond between the two atoms). Thus, the classical contribution to the specific heat from vibrational degrees of freedom is

$$(c_V)_{\text{vibration}} = (n - 1)R. \quad (5.407)$$

So, do any of the rotational and vibrational degrees of freedom actually make a contribution to the specific heats of gases at room temperature, once quantum effects have been taken into consideration? We can answer this question by examining just one piece of data. Figure 5.4 shows the infrared absorption spectrum of hydrogen chloride gas. The absorption lines correspond to simultaneous transitions between different vibrational and rotational energy levels. Hence, this is usually called a *vibration-rotation spectrum*. The missing line at about 3.47 microns corresponds to a pure vibrational transition from the ground state to the first excited state. (Pure vibrational transitions are forbidden; hydrogen chloride molecules always have to simultaneously change their rotational

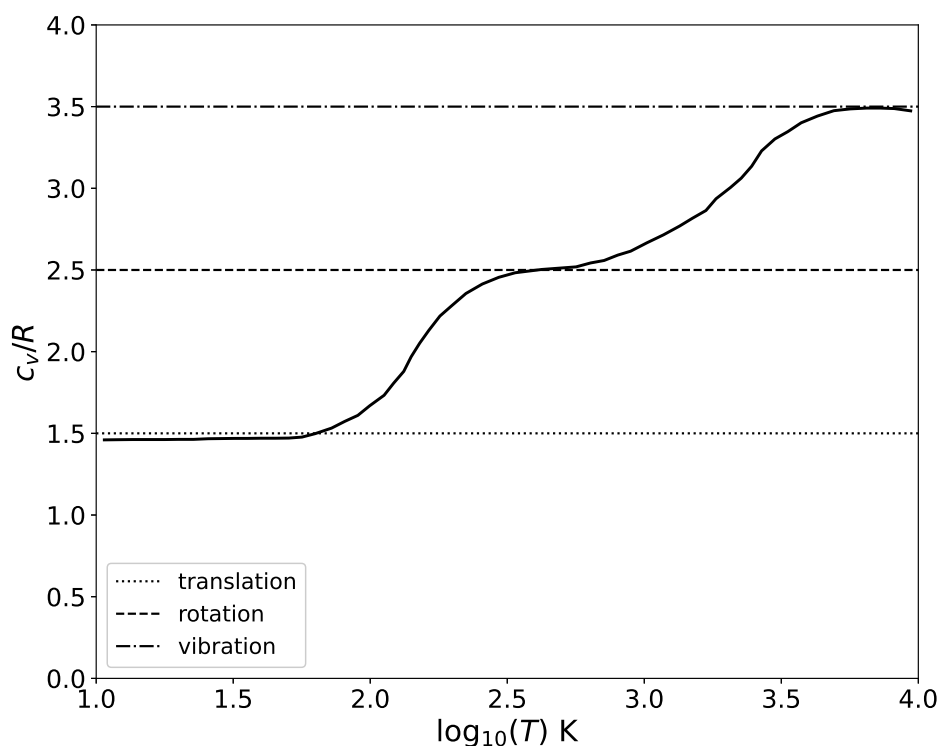


Figure 5.5: The molar heat capacity at constant volume of gaseous molecular hydrogen versus temperature.

energy level if they are to couple effectively to electromagnetic radiation.) The longer wavelength absorption lines correspond to vibrational transitions in which there is a simultaneous decrease in the rotational energy level. Likewise, the shorter wavelength absorption lines correspond to vibrational transitions in which there is a simultaneous increase in the rotational energy level. It is clear that the rotational energy levels are more closely spaced than the vibrational energy levels. The pure vibrational transition gives rise to absorption at about 3.47 microns, which corresponds to infrared radiation of frequency 8.5×10^{13} hertz with an associated radiation “temperature” of 4,100K. We conclude that the vibrational degrees of freedom of hydrogen chloride, or any other small molecule, are frozen out at room temperature. The rotational transitions split the vibrational lines by about 0.2 microns. This implies that pure rotational transitions would be associated with infrared radiation of frequency 5×10^{12} hertz and corresponding radiation “temperature” 240K. We conclude that the rotational degrees of freedom of hydrogen chloride, or any other small molecule, are not frozen out at room temperature, and probably contribute the classical $(1/2)R$ to the molar specific heat. There is one proviso, however. Linear molecules (like hydrogen chloride) effectively only have two rotational degrees of freedom (instead of the usual three), because of the very small moment of inertia of such molecules along the line of centers of the atoms.

Figure 5.5 shows the variation of the molar heat capacity at constant volume of gaseous molecular hydrogen (i.e., H_2) with temperature. The expected contribution from the translational degrees of freedom is $(3/2)R$ (there are three translational degrees of freedom per molecule). The expected contribution at high temperatures from the rotational degrees of freedom is R (there are effectively two rotational degrees of freedom per molecule). Finally, the expected contribution at high temperatures from the vibrational degrees of freedom is R (there is one vibrational degree of freedom per molecule). It can be seen that, as the temperature rises, the rotational, and then the vibrational, degrees of freedom eventually make their full classical contributions to the heat capacity.

5.5.9 Maxwell Velocity Distribution

Consider a molecule of mass m in a gas that is sufficiently dilute for the intermolecular forces to be negligible (i.e., an ideal gas). The energy of the molecule is written

$$\epsilon = \frac{\mathbf{p}^2}{2m} + \epsilon^{\text{int}}, \quad (5.408)$$

where \mathbf{p} is its momentum vector, and ϵ^{int} is its internal (i.e., non-translational) energy. The latter energy is due to molecular rotation, vibration, et cetera. Translational degrees of freedom can be treated classically to an excellent approximation, whereas internal degrees of freedom usually require a quantum-mechanical approach. Classically, the probability of finding the molecule in a given internal state with a position vector in the range \mathbf{r} to $\mathbf{r} + d\mathbf{r}$, and a momentum vector in the range \mathbf{p} to $\mathbf{p} + d\mathbf{p}$, is proportional to the number of cells (of “volume” h_0) contained in the corresponding region of phase-space, weighted by the Boltzmann factor. (See Section 5.5.5.) In fact, because classical phase-space is divided up into uniform cells, the number of cells is just proportional to the “volume” of the region under consideration. (See Section 5.4.4.) This “volume” is written $d^3\mathbf{r} d^3\mathbf{p}$. Thus, the probability of finding the molecule in a given internal state s is

$$P_s(\mathbf{r}, \mathbf{p}) d^3\mathbf{r} d^3\mathbf{p} \propto \exp\left(-\frac{p^2}{2mk_B T}\right) \exp\left(-\frac{\epsilon_s^{\text{int}}}{k_B T}\right) d^3\mathbf{r} d^3\mathbf{p}, \quad (5.409)$$

where P_s is a probability density defined in the usual manner. The probability $P(\mathbf{r}, \mathbf{p}) d^3\mathbf{r} d^3\mathbf{p}$ of finding the molecule in any internal state with position and momentum vectors in the specified range is obtained by summing the previous expression over all possible internal states. The sum over $\exp(\epsilon_s^{\text{int}}/k_B T)$ just contributes a constant of proportionality (because the internal states do not depend on \mathbf{r} or \mathbf{p}), so

$$P(\mathbf{r}, \mathbf{p}) d^3\mathbf{r} d^3\mathbf{p} \propto \exp\left(-\frac{p^2}{2mk_B T}\right) d^3\mathbf{r} d^3\mathbf{p}. \quad (5.410)$$

Of course, we can multiply this probability by the total number of molecules, N , in order to obtain the mean number of molecules with position and momentum vectors in the specified range.

Suppose that we now wish to determine $f(\mathbf{r}, \mathbf{v}) d^3\mathbf{r} d^3\mathbf{v}$; that is, the mean number of molecules with positions between \mathbf{r} and $\mathbf{r} + d\mathbf{r}$, and velocities in the range \mathbf{v} and $\mathbf{v} + d\mathbf{v}$. Because $\mathbf{v} = \mathbf{p}/m$, it is easily seen that

$$f(\mathbf{r}, \mathbf{v}) d^3\mathbf{r} d^3\mathbf{v} = C \exp\left(-\frac{mv^2}{2k_B T}\right) d^3\mathbf{r} d^3\mathbf{v}, \quad (5.411)$$

where C is a constant of proportionality. This constant can be determined by the condition

$$\int_{(\mathbf{r})} \int_{(\mathbf{v})} f(\mathbf{r}, \mathbf{v}) d^3\mathbf{r} d^3\mathbf{v} = N. \quad (5.412)$$

In other words, the sum over molecules with all possible positions and velocities must give the total number of molecules, N . The integral over the molecular position coordinates just gives the volume, V , of the gas, because the Boltzmann factor is independent of position. The integration over the velocity coordinates can be reduced to the product of three identical integrals (one for v_x , one for v_y , and one for v_z), so we have

$$C V \left[\int_{-\infty}^{\infty} \exp\left(-\frac{m v_z^2}{2 k_B T}\right) dv_z \right]^3 = N. \quad (5.413)$$

Now,

$$\int_{-\infty}^{\infty} \exp\left(-\frac{m v_z^2}{2 k_B T}\right) dv_z = \sqrt{\frac{2 k_B T}{m}} \int_{-\infty}^{\infty} \exp(-y^2) dy = \sqrt{\frac{2\pi k_B T}{m}}, \quad (5.414)$$

so $C = (N/V) (m/2\pi k_B T)^{3/2}$. Thus, the properly normalized distribution function for molecular velocities is written

$$f(\mathbf{v}) d^3\mathbf{r} d^3\mathbf{v} = n \left(\frac{m}{2\pi k_B T} \right)^{3/2} \exp\left(-\frac{m v^2}{2 k_B T}\right) d^3\mathbf{r} d^3\mathbf{v}. \quad (5.415)$$

Here, $n = N/V$ is the number density of the molecules. We have omitted the variable \mathbf{r} in the argument of f , because f clearly does not depend on position. In other words, the distribution of molecular velocities is uniform in space. This is hardly surprising, because there is nothing to distinguish one region of space from another in our calculation. The previous distribution is called the *Maxwell velocity distribution*, because it was discovered by James Clark Maxwell in the middle of the nineteenth century. The average number of molecules per unit volume with velocities in the range \mathbf{v} to $\mathbf{v} + d\mathbf{v}$ is obviously $f(\mathbf{v}) d^3\mathbf{v}$. Note that $\int f(\mathbf{v}) d^3\mathbf{v} = n$.

Let us consider the distribution of a given component of velocity; the z -component (say). Suppose that $g(v_z) dv_z$ is the average number of molecules per unit volume with the z -component of velocity in the range v_z to $v_z + dv_z$, irrespective of the values of their other velocity components. It is fairly obvious that this distribution is obtained from the Maxwell distribution by summing (integrating actually) over all possible values of v_x and v_y , with v_z in the specified range. Thus,

$$g(v_z) dv_z = \int_{(v_x)} \int_{(v_y)} f(\mathbf{v}) d^3\mathbf{v}. \quad (5.416)$$

This gives

$$\begin{aligned} g(v_z) dv_z &= n \left(\frac{m}{2\pi k_B T} \right)^{3/2} \int_{(v_x)} \int_{(v_y)} \exp\left[-\left(\frac{m}{2 k_B T}\right) (v_x^2 + v_y^2 + v_z^2)\right] dv_x dv_y dv_z \\ &= n \left(\frac{m}{2\pi k_B T} \right)^{3/2} \exp\left(-\frac{m v_z^2}{2 k_B T}\right) \left[\int_{-\infty}^{\infty} \exp\left(-\frac{m v_x^2}{2 k_B T}\right) \right]^2 \end{aligned}$$

$$= n \left(\frac{m}{2\pi k_B T} \right)^{3/2} \exp \left(-\frac{m v_z^2}{2 k_B T} \right) \left(\sqrt{\frac{2\pi k_B T}{m}} \right)^2, \quad (5.417)$$

or

$$g(v_z) dv_z = n \left(\frac{m}{2\pi k_B T} \right)^{1/2} \exp \left(-\frac{m v_z^2}{2 k_B T} \right) dv_z. \quad (5.418)$$

Of course, this expression is properly normalized, so that

$$\int_{-\infty}^{\infty} g(v_z) dv_z = n. \quad (5.419)$$

It is clear that each component (because there is nothing special about the z -component) of the velocity is distributed with a Gaussian probability distribution (see Section 5.1.7), centered on a mean value

$$\langle v_z \rangle = 0, \quad (5.420)$$

with variance

$$\langle v_z^2 \rangle = \frac{k_B T}{m}. \quad (5.421)$$

Equation (5.420) implies that each molecule is just as likely to be moving in the plus z -direction as in the minus z -direction. Equation (5.421) can be rearranged to give

$$\left\langle \frac{1}{2} m v_z^2 \right\rangle = \frac{1}{2} k_B T, \quad (5.422)$$

in accordance with the equipartition theorem. (See Section 5.5.5.)

Note that Equation (5.415) can be rewritten

$$\frac{f(\mathbf{v}) d^3 \mathbf{v}}{n} = \left[\frac{g(v_x) dv_x}{n} \right] \left[\frac{g(v_y) dv_y}{n} \right] \left[\frac{g(v_z) dv_z}{n} \right], \quad (5.423)$$

where $g(v_x)$ and $g(v_y)$ are defined in an analogous way to $g(v_z)$. Thus, the probability that the velocity lies in the range \mathbf{v} to $\mathbf{v} + d\mathbf{v}$ is just equal to the product of the probabilities that the velocity components lie in their respective ranges. In other words, the individual velocity components act like statistically independent variables.

Suppose that we now wish to calculate $F(v) dv$; that is, the average number of molecules per unit volume with a speed $v = |\mathbf{v}|$ in the range v to $v + dv$. It is obvious that we can obtain this quantity by summing over all molecules with speeds in this range, irrespective of the direction of their velocities. Thus,

$$F(v) dv = \int f(\mathbf{v}) d^3 \mathbf{v}, \quad (5.424)$$

where the integral extends over all velocities satisfying

$$v < |\mathbf{v}| < v + dv. \quad (5.425)$$

This inequality is satisfied by a spherical shell of radius v and thickness dv in velocity space. Because $f(\mathbf{v})$ only depends on $|v|$, so $f(\mathbf{v}) \equiv f(v)$, the previous integral is just $f(v)$ multiplied by the volume of the spherical shell in velocity space. So,

$$F(v) dv = f(v) 4\pi v^2 dv, \quad (5.426)$$

which gives

$$F(v) dv = 4\pi n \left(\frac{m}{2\pi k_B T} \right)^{3/2} v^2 \exp \left(-\frac{m v^2}{2 k_B T} \right) dv. \quad (5.427)$$

This result is known as *Maxwell's distribution of molecular speeds*. Of course, it is properly normalized, so that

$$\int_0^\infty F(v) dv = n. \quad (5.428)$$

Note that the Maxwell distribution exhibits a maximum at some non-zero value of v . The reason for this is quite simple. As v increases, the Boltzmann factor decreases, but the volume of phase-space available to the molecule (which is proportional to v^2) increases; the net result is a distribution with a non-zero maximum.

The mean molecular speed is given by

$$\langle v \rangle = \frac{1}{n} \int_0^\infty F(v) v dv. \quad (5.429)$$

Thus, we obtain

$$\langle v \rangle = 4\pi \left(\frac{m}{2\pi k_B T} \right)^{3/2} \int_0^\infty v^3 \exp \left(-\frac{m v^2}{2 k_B T} \right) dv, \quad (5.430)$$

or

$$\langle v \rangle = 4\pi \left(\frac{m}{2\pi k_B T} \right)^{3/2} \left(\frac{2 k_B T}{m} \right)^2 \int_0^\infty y^3 \exp(-y^2) dy. \quad (5.431)$$

Now,

$$\int_0^\infty y^3 \exp(-y^2) dy = \frac{1}{2}, \quad (5.432)$$

so

$$\langle v \rangle = \sqrt{\frac{8 k_B T}{\pi m}}. \quad (5.433)$$

A similar calculation gives

$$v_{\text{rms}} = [\langle v^2 \rangle]^{1/2} = \sqrt{\frac{3 k_B T}{m}}. \quad (5.434)$$

However, this result can also be obtained from the equipartition theorem. (See Section 5.5.5.) Because

$$\left\langle \frac{1}{2} m v^2 \right\rangle = \left\langle \frac{1}{2} \frac{p_x^2}{m} \right\rangle + \left\langle \frac{1}{2} \frac{p_y^2}{m} \right\rangle + \left\langle \frac{1}{2} \frac{p_z^2}{m} \right\rangle = 3 \left(\frac{1}{2} k_B T \right), \quad (5.435)$$

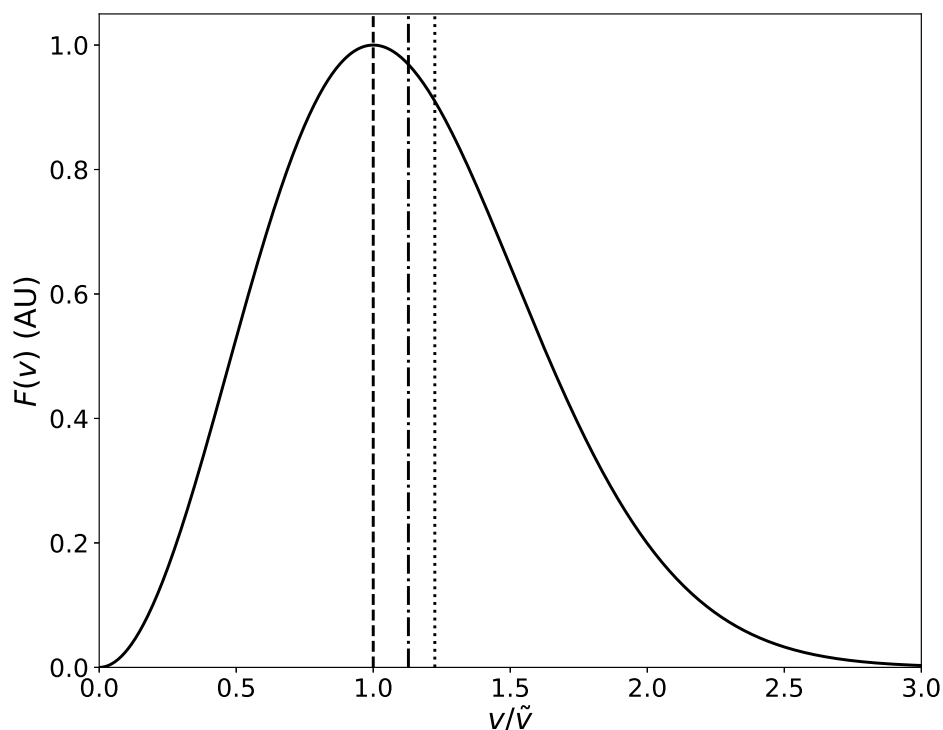


Figure 5.6: The Maxwell velocity distribution as a function of molecular speed, in units of the most probable speed (\tilde{v}). The dashed, dash-dotted, and dotted lines indicates the most probable speed, the mean speed, and the root-mean-square speed, respectively.

then Equation (5.434) follows immediately. It is easily demonstrated that the most probable molecular speed (i.e., the maximum of the Maxwell distribution function) is

$$\tilde{v} = \sqrt{\frac{2k_B T}{m}}. \quad (5.436)$$

The speed of sound in an ideal gas is given by

$$v_s = \sqrt{\frac{\gamma p}{\rho}}, \quad (5.437)$$

where γ is the ratio of specific heats. (See Section 5.2.9.) This can also be written

$$v_s = \sqrt{\frac{\gamma k_B T}{m}}, \quad (5.438)$$

because $p = nk_B T$ and $\rho = nm$. It is clear that the various average speeds that we have just calculated are all of order the sound speed (i.e., a few hundred meters per second at room temperature).

In ordinary air ($\gamma = 1.4$) the sound speed is about 84% of the most probable molecular speed, and about 74% of the mean molecular speed. Because sound waves ultimately propagate via molecular motion, it makes sense that they travel at slightly less than the most probable and mean molecular speeds.

Figure 5.6 shows the Maxwell velocity distribution as a function of molecular speed in units of the most probable speed. Also shown are the mean speed and the root-mean-square speed.

5.6 Standing-Wave States

5.6.1 Counting Standing-Wave States

Consider a three-dimensional standing wave confined in a cubic box that extends from $x = 0$ to $x = a$, from $y = 0$ to $y = a$, and from $z = 0$ to $z = a$. (See Section 4.4.2.) The wavefunction, $\psi(x, y, z)$, must satisfy the boundary conditions

$$\psi(0, y, z) = \psi(a, y, z) = 0, \quad (5.439)$$

$$\psi(x, 0, z) = \psi(x, a, z) = 0, \quad (5.440)$$

$$\psi(x, y, 0) = \psi(x, y, a) = 0. \quad (5.441)$$

Thus, standing-wave solutions of the form

$$\psi(x, y, z) = \psi_0 \sin(k_x x) \sin(k_y y) \sin(k_z z) \quad (5.442)$$

are only acceptable if

$$k_x = n_x \frac{\pi}{a}, \quad (5.443)$$

$$k_y = n_y \frac{\pi}{a}, \quad (5.444)$$

$$k_z = n_z \frac{\pi}{a}, \quad (5.445)$$

where n_x , n_y , and n_z are positive integers. (Note that negative values of n_x do not give rise to wave states that are physically distinct from the corresponding positive values, et cetera.) It follows that k_x , k_y , and k_z are all quantized in units of π/a .

Now,

$$\frac{\partial n_x}{\partial k_x} = \frac{a}{\pi}, \quad (5.446)$$

$$\frac{\partial n_y}{\partial k_y} = \frac{a}{\pi}, \quad (5.447)$$

$$\frac{\partial n_z}{\partial k_z} = \frac{a}{\pi}. \quad (5.448)$$

Thus, the number of translational wave states that are such that k_x lies between k_x and $k_x + dk_x$, k_y lies between k_y and $k_y + dk_y$, and k_z lies between k_z and $k_z + dk_z$, is

$$N(\mathbf{k}) dk_x dk_y dk_z = \left(\frac{\partial n_x}{\partial k_x} dk_x \right) \left(\frac{\partial n_y}{\partial k_y} dk_y \right) \left(\frac{\partial n_z}{\partial k_z} dk_z \right) = \left(\frac{a}{\pi} \right)^3 dk_x dk_y dk_z. \quad (5.449)$$

Note that

$$N(\mathbf{k}) = \left(\frac{a}{\pi} \right)^3 \quad (5.450)$$

is independent of the wavevector $\mathbf{k} = (k_x, k_y, k_z)$, because the allowed wave states are uniformly distributed in \mathbf{k} -space.

The wavenumber is defined $k = |\mathbf{k}|$. The number of translational wave states such that k lies between k and $k + dk$ is denoted $\rho(k) dk$, where $\rho(k)$ is termed the *density of states*. Now, $\rho(k) dk$ is the number of wave states that lie in an octant of a spherical annulus in \mathbf{k} -space whose inner radius is k , and whose outer radius is $k + dk$. We have to take an octant of the annulus because only wave states characterized by positive values of n_x , n_y , and n_z have physical significance. (See Section 4.4.3.) The volume of the octant in \mathbf{k} -space is

$$\mathcal{V} = \frac{1}{8} 4\pi k^2 dk. \quad (5.451)$$

Hence,

$$\rho(k) dk = N(\mathbf{k}) \mathcal{V} = \left(\frac{a}{\pi} \right)^3 \frac{1}{8} 4\pi k^2 dk, \quad (5.452)$$

which implies that

$$\rho(k) = \frac{V k^2}{2\pi^2}, \quad (5.453)$$

where $V = a^3$ is the volume of the box. Although the previous expression was derived for the special case of a cubic box, we shall assume that it is valid for a macroscopic box of any shape. This assumption is reasonable provided that the wavelengths of most of the standing waves confined in the box are much smaller than the dimensions of the box (i.e., provided that n_x , n_y , and n_z are all typically much greater than unity).

Consider electromagnetic waves confined in a box. Such waves satisfy the dispersion relation

$$\omega = k c, \quad (5.454)$$

where c is the speed of light in vacuum. (See Section 2.4.4.) Note that c is not a function of ω . Let $\rho(\omega) d\omega$ be the number of translational electromagnetic wave states for which ω lies between ω and $\omega + d\omega$. It follows that

$$\rho(\omega) d\omega = 2\rho(k) \frac{dk}{d\omega} d\omega = 2\rho(k) \frac{1}{c} d\omega, \quad (5.455)$$

which yields

$$\rho(\omega) = 2 \frac{\rho(k)}{c} = 2 \frac{V k^2}{2\pi^2 c}, \quad (5.456)$$

giving

$$\rho(\omega) = \frac{V}{\pi^2} \frac{\omega^2}{c^3}, \quad (5.457)$$

where use has been made of Equations (5.453) and (5.454). Here, the factor of 2 in Equation (5.455) is required because electromagnetic waves are transverse waves, so there are two independent polarization states for each allowed translational state. (See Section 2.4.4.)

Consider sound waves propagating through a solid. Such waves satisfy the dispersion relation

$$\omega = v_s k, \quad (5.458)$$

where v_s is the sound speed. Note that v_s is not (usually) a function of ω . However, solids support both transverse and longitudinal sound waves (unlike gases, which only support longitudinal waves). Of course, for transverse waves, there are two independent polarization states for each allowed translational state. However, for longitudinal waves, there is only one polarization state for each allowed translational state. Thus, by analogy with electromagnetic waves, the density of transverse sound wave states is

$$\rho_t(\omega) = \frac{V}{\pi^2} \frac{\omega^2}{v_{st}^3}, \quad (5.459)$$

where v_{ts} is the characteristic phase velocity of transverse waves, whereas the density of longitudinal sound wave states is

$$\rho_l(\omega) = \frac{V}{2\pi^2} \frac{\omega^2}{v_{sl}^3}, \quad (5.460)$$

where v_{ls} is the characteristic phase velocity of longitudinal waves. The total density of sound wave states, irrespective of the wave polarization, is

$$\rho(\omega) = \frac{3V}{2\pi^2} \frac{\omega^2}{v_s^3}, \quad (5.461)$$

where

$$\frac{1}{v_s^3} = \frac{2}{3} \frac{1}{v_{st}^3} + \frac{1}{3} \frac{1}{v_{sl}^3}. \quad (5.462)$$

Here, v_s is the average sound speed.

Finally, consider electrons of mass m_e confined in a box. According to quantum mechanics, electrons have wavelike properties such that the electron energy, ϵ , is related to the wavenumber, k , according to the dispersion relation

$$\epsilon = \frac{\hbar^2 k^2}{2m_e}. \quad (5.463)$$

(See Section 4.4.2.) Let $\rho(\epsilon) d\epsilon$ be the number of translational electron states for which ϵ lies between ϵ and $\epsilon + d\epsilon$. It follows that

$$\rho(\epsilon) d\epsilon = \rho(k) \frac{dk}{d\epsilon} d\epsilon = \rho(k) \frac{(2m_e)^{1/2}}{2\hbar \epsilon^{1/2}} d\epsilon. \quad (5.464)$$

However, according to the *Pauli exclusion principle* (see Section 4.4.3), only two electrons (corresponding to a spin-up electron, and a spin-down electron) can be put into each translational state.

Hence, reinterpreting $\rho(\epsilon) d\epsilon$ as the number of electrons whose energies lies between ϵ and $\epsilon + d\epsilon$, we get

$$\rho(\epsilon) = 2 \frac{V k^2}{2\pi^2} \frac{(2 m_e)^{1/2}}{2 \hbar \epsilon^{1/2}}, \quad (5.465)$$

which gives

$$\rho(\epsilon) = \frac{\sqrt{2} V m_e^{3/2} \epsilon^{1/2}}{\pi^2 \hbar^3}, \quad (5.466)$$

where use has been made of Equations (5.453) and (5.463).

5.6.2 Planck Radiation Law

Consider electromagnetic radiation inside a box whose walls are held at the constant temperature T . We know that electromagnetic radiation of angular frequency ω is quantized into photons whose energy is $\epsilon = \hbar \omega$. (See Section 3.3.8 and 4.1.2.) Thus, given that photons are indivisible, the allowed energy levels of such radiation are equally spaced, with spacing $\hbar \omega$. In this respect, each frequency state acts like a harmonic oscillator of angular frequency ω . (See Section 4.3.7.) According to Equation (5.390), the mean energy of a harmonic oscillator of angular frequency ω that is in thermal equilibrium with a heat reservoir of temperature T is

$$\langle \epsilon \rangle = \frac{\hbar \omega}{\exp(\hbar \omega / k_B T) - 1}. \quad (5.467)$$

Here, we have neglected the zero-point energy, $(1/2) \hbar \omega$, in Equation (5.390) because there is no electromagnetic zero-point energy.

Let $u(\omega)$ be the electromagnetic energy per unit volume associated with electromagnetic waves whose angular frequencies lie between ω and $\omega + d\omega$. It follows that

$$V u(\omega) d\omega = \rho(\omega) \langle \epsilon \rangle d\omega, \quad (5.468)$$

where $\rho(\omega) d\omega$ is the number of electromagnetic wave states whose angular frequencies lie between ω and $\omega + d\omega$. Making use of Equations (5.457) and (5.467), we deduce that

$$u(\omega) = \frac{\hbar \omega^3}{\pi^2 c^3 [\exp(\hbar \omega / k_B T) - 1]}. \quad (5.469)$$

This result is known as the *Planck radiation law*, after Max Planck who first obtained it in 1900.

Consider the classical limit $\hbar \rightarrow 0$. In this limit, the previous expression becomes

$$u(\omega) = \frac{k_B T \omega^2}{\pi^2 c^3}. \quad (5.470)$$

This result is known as the *Rayleigh-Jeans radiation law*, after Lord Rayleigh and James Jeans who derived it in the first decade of the twentieth century. The Rayleigh-Jeans law is equivalent

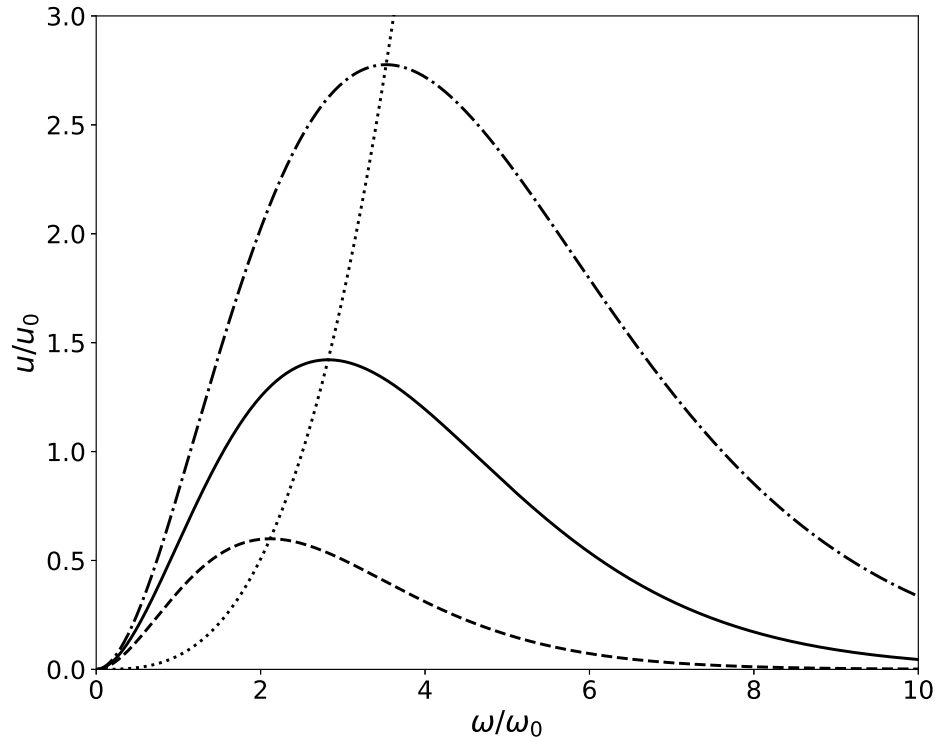


Figure 5.7: The Planck radiation law. Here, $\omega_0 = k_B T_0 / \hbar$ and $u_0 = \hbar \omega_0^3 / (\pi^2 c^3)$, where T_0 is an arbitrary scale temperature. The dashed, solid, and dash-dotted curves show u/u_0 for $T/T_0 = 0.75$, 1.0, and 1.25, respectively. The dotted curve shows the locus of the peak emission frequency.

to the assumption that each electromagnetic wave state possesses the classical energy $k_B T$ predicted by the equipartition theorem. (See Section 5.5.5.) The total classical energy density of electromagnetic radiation is given by

$$u_{\text{tot}} = \int_0^{\infty} u(\omega) d\omega = \frac{k_B T}{\pi^2 c^3} \int_0^{\infty} \omega^2 d\omega. \quad (5.471)$$

This is an integral that obviously does not converge. Thus, according to classical physics, the total energy density of electromagnetic radiation inside an enclosed cavity is infinite. This is clearly an absurd result. In fact, this prediction is known as the *ultra-violet catastrophe*, because the Rayleigh-Jeans law usually starts to diverge badly from experimental observations (by overestimating the amount of radiation) in the ultra-violet region of the spectrum.

The Planck radiation law approximates to the classical Rayleigh-Jeans law for $\hbar \omega \ll k_B T$, peaks at about $\hbar \omega \simeq 3 k_B T$, and falls off exponentially for $\hbar \omega \gg k_B T$. See Figure 5.7. The exponential fall-off at high frequencies ensures that the total energy density of electromagnetic radiation inside an enclosed cavity remains finite. The reason for the fall-off that it is very difficult

for a thermal fluctuation to create a photon with an energy greatly in excess of $k_B T$, because $k_B T$ is the characteristic energy associated with such fluctuations.

5.6.3 Black-Body Radiation

Suppose that we were to make a small hole in the wall of the enclosure described in the previous section, and were then to observe the emitted radiation. A small hole is the best approximation in physics to a *black-body*, which is defined as an object that absorbs, and, therefore, emits, radiation perfectly at all wavelengths. What is the power radiated by the hole?

The power density inside the enclosure can be written

$$u(\omega) d\omega = \hbar \omega n(\omega) d\omega, \quad (5.472)$$

where $n(\omega)$ is the mean number of photons per unit volume whose frequencies lie in the range ω to $\omega + d\omega$. The radiation field inside the enclosure is isotropic (we are assuming that the hole is sufficiently small that it does not distort the field). It follows that the mean number of photons per unit volume whose frequencies lie in the specified range, and whose directions of propagation subtend an angle in the range θ to $\theta + d\theta$ with the normal to the hole, is

$$n(\omega, \theta) d\omega d\theta = n(\omega) d\omega g(\theta) d\theta, \quad (5.473)$$

where $g(\theta) d\theta = (1/2) \sin \theta d\theta$ is the fractional range of solid angle in the specified range of directions. (See Section 5.3.2.) The previous two equations give

$$\hbar \omega n(\omega, \theta) = \frac{1}{2} u(\omega) \sin \theta. \quad (5.474)$$

Photons travel at the speed of light, so the power per unit area escaping from the hole in the frequency range ω to $\omega + d\omega$ is

$$P(\omega) d\omega = \int_0^{\pi/2} c \cos \theta \hbar \omega n(\omega, \theta) d\omega d\theta, \quad (5.475)$$

where $c \cos \theta$ is the component of the photon velocity in the direction of the hole. This gives

$$P(\omega) d\omega = c u(\omega) d\omega \frac{1}{2} \int_0^{\pi/2} \cos \theta \sin \theta d\theta = \frac{c}{4} u(\omega) d\omega, \quad (5.476)$$

so

$$P(\omega) d\omega = \frac{\hbar}{4\pi^2 c^2} \frac{\omega^3 d\omega}{\exp(\hbar \omega/k_B T) - 1} \quad (5.477)$$

is the power per unit area radiated by a black-body in the frequency range ω to $\omega + d\omega$.

A black-body is very much an idealization. The power spectra of real radiating bodies can deviate quite substantially from black-body spectra. Nevertheless, we can make some useful predictions using this model. The black-body power spectrum peaks when $\hbar \omega \simeq 3 k_B T$, implying that the peak radiation frequency scales linearly with the temperature of the body. In other words,

hot bodies tend to radiate at higher frequencies than cold bodies. This result (in particular, the linear scaling) is known as *Wien's displacement law*, after Wilhelm Wien who derived it in 1893, and allows us to estimate the surface temperatures of stars from their colors (surprisingly enough, stars are fairly good black-bodies). Table 5.2 shows some stellar temperatures determined by this method (in fact, the whole emission spectrum is fitted to a black-body spectrum). It can be seen that the apparent colors (which correspond quite well to the colors of the peak radiation) scan the whole visible spectrum, from red to blue, as the stellar surface temperatures gradually rise.

Name	Constellation	Surface Temp. (K)	Color
Antares	Scorpio	3300	Very Red
Aldebaran	Taurus	3800	Reddish Yellow
Sun		5770	Yellow
Procyon	Canis Minor	6570	Yellowish White
Sirius	Canis Major	9250	White
Rigel	Orion	11,200	Bluish White

Table 5.2: Physical properties of some well-known stars.

Probably the most famous black-body spectrum is cosmological in origin. Just after the “big bang,” the universe was essentially a “fireball,” with the energy associated with radiation completely dominating that associated with matter. The early universe was also fairly well described by equilibrium statistical thermodynamics, which means that the radiation had a black-body spectrum. As the universe expanded, the radiation was gradually Doppler shifted to ever larger wavelengths (in other words, the radiation did work against the expansion of the universe, and, thereby, lost energy, but its spectrum remained invariant). Nowadays, this primordial radiation is detectable as a faint microwave background that pervades the whole universe. The cosmic microwave background was discovered accidentally by Arno Penzias and Robert Wilson in 1964. For many years, it was difficult to measure the full spectrum of the microwave background with any degree of precision, because of strong absorption and scattering of microwaves by the Earth’s atmosphere. However, all of this changed when the COBE satellite was launched in 1989. It took precisely nine minutes to measure the perfect black-body spectrum reproduced in Figure 5.8. The data shown in the figure can be fitted to a black-body curve of characteristic temperature 2.735K. In a very real sense, this can be regarded as the “temperature of the universe.”

5.6.4 Stefan-Boltzmann Law

The total power radiated per unit area by a black-body at all frequencies is given by

$$P_{\text{tot}}(T) = \int_0^{\infty} P(\omega) d\omega = \frac{\hbar}{4\pi^2 c^2} \int_0^{\infty} \frac{\omega^3 d\omega}{\exp(\hbar \omega/k_B T) - 1}, \quad (5.478)$$

or

$$P_{\text{tot}}(T) = \frac{k_B^4 T^4}{4\pi^2 c^2 \hbar^3} \int_0^{\infty} \frac{\eta^3 d\eta}{e^\eta - 1}, \quad (5.479)$$

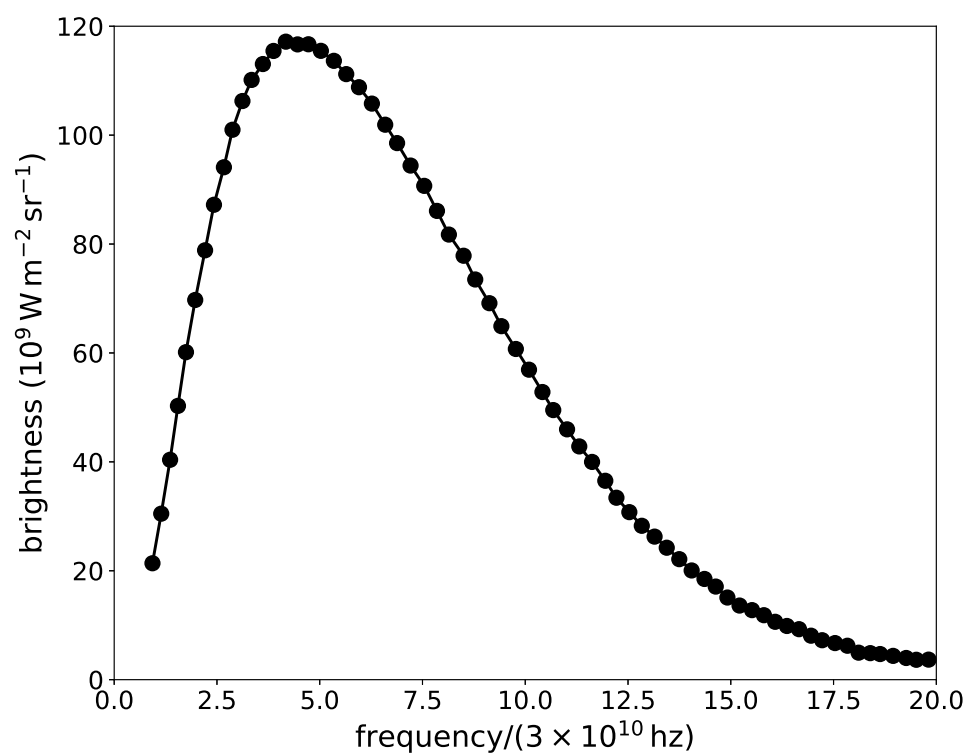


Figure 5.8: Cosmic background radiation spectrum measured by the Far Infrared Absolute Spectrometer (FIRAS) aboard the Cosmic Background Explorer satellite (COBE). The fit is to a black-body spectrum of characteristic temperature $2.735 \pm 0.06 \text{ K}$.

where $\eta = \hbar \omega / k_B T$. The previous integral can be looked up in standard reference books on integrals. In fact,

$$\int_0^\infty \frac{\eta^3 d\eta}{e^\eta - 1} = \frac{\pi^4}{15}. \quad (5.480)$$

Thus, the total power radiated per unit area by a black-body is

$$P_{\text{tot}}(T) = \frac{\pi^2}{60} \frac{k_B^4}{c^2 \hbar^3} T^4 = \sigma T^4. \quad (5.481)$$

This T^4 dependence of the radiated power is called the *Stefan-Boltzmann law*, after Josef Stefan, who first obtained it experimentally 1877, and Ludwig Boltzmann, who first derived it theoretically in 1884. The parameter

$$\sigma = \frac{\pi^2}{60} \frac{k_B^4}{c^2 \hbar^3} = 5.67 \times 10^{-8} \text{ W m}^{-2} \text{ K}^{-4}, \quad (5.482)$$

is known as the *Stefan-Boltzmann constant*.

We can use the Stefan-Boltzmann law to estimate the temperature of the Earth from first principles. The Sun is a ball of glowing gas of radius $R_\odot \simeq 7 \times 10^5$ km and surface temperature $T_\odot \simeq 5770$ K. Its luminosity is

$$L_\odot = 4\pi R_\odot^2 \sigma T_\odot^4, \quad (5.483)$$

according to the Stefan-Boltzmann law. The Earth is a globe of radius $R_\oplus \simeq 6000$ km located an average distance $r_\oplus \simeq 1.5 \times 10^8$ km from the Sun. The Earth intercepts an amount of energy

$$P_\oplus = L_\odot \frac{\pi R_\oplus^2 / r_\oplus^2}{4\pi} \quad (5.484)$$

per second from the Sun's radiative output; that is, the power output of the Sun reduced by the ratio of the solid angle subtended by the Earth at the Sun to the total solid angle 4π . The Earth absorbs this energy, and then re-radiates it at longer wavelengths. The luminosity of the Earth is

$$L_\oplus = 4\pi R_\oplus^2 \sigma T_\oplus^4, \quad (5.485)$$

according to the Stefan-Boltzmann law, where T_\oplus is the average temperature of the Earth's surface. Here, we are ignoring any surface temperature variations between polar and equatorial regions, or between day and night. In a steady state, the luminosity of the Earth must balance the radiative power input from the Sun, so, equating L_\oplus and P_\oplus , we arrive at

$$T_\oplus = \left(\frac{R_\odot}{2r_\oplus} \right)^{1/2} T_\odot. \quad (5.486)$$

Remarkably, the ratio of the Earth's surface temperature to that of the Sun depends only on the Earth-Sun distance and the solar radius. The previous expression yields $T_\oplus \simeq 279$ K or 6° C (or 43° F). This is slightly on the cold side, by a few degrees, because of the greenhouse action of the Earth's atmosphere, which was neglected in our calculation. Nevertheless, it is quite encouraging that such a crude calculation comes so close to the correct answer.

5.6.5 Specific Heats of Solids

Consider a simple solid containing N atoms. Now, atoms in solids cannot translate (unlike those in gases), but are free to vibrate about their equilibrium positions. Such vibrations are termed *lattice vibrations*, and can be thought of as sound waves propagating through the crystal lattice. Each atom is specified by three independent position coordinates, and three corresponding momentum coordinates. Let us only consider small-amplitude vibrations. In this case, we can expand the potential energy of interaction between the atoms to give an expression that is quadratic in the atomic displacements from their equilibrium positions. It is always possible to perform a *normal mode analysis* of the oscillations. In effect, we can find $3N$ independent modes of oscillation of the solid. Each mode has its own particular oscillation frequency, and its own particular pattern of atomic displacements. Any general oscillation can be written as a linear combination of these *normal modes*. Let q_i be the (appropriately normalized) amplitude of the i th normal mode, and p_i the corresponding momentum. In *normal-mode coordinates*, the internal energy of the lattice vibrations takes the particularly simple form

$$U = \frac{1}{2} \sum_{i=1,3N} (p_i^2 + \omega_i^2 q_i^2), \quad (5.487)$$

where ω_i is the (angular) oscillation frequency of the i th normal mode. It is clear that, when expressed in normal-mode coordinates, the linearized lattice vibrations are equivalent to $3N$ independent harmonic oscillators. (Of course, each oscillator corresponds to a different normal mode.)

The typical value of ω_i is the (angular) frequency of a sound wave propagating through the lattice. Sound wave frequencies are far lower than the typical vibration frequencies of gaseous molecules. In the latter case, the mass involved in the vibration is simply that of the molecule, whereas in the former case the mass involved is that of very many atoms (because lattice vibrations are non-localized). The strength of interatomic bonds in gaseous molecules is similar to those in solids, so we can use the estimate $\omega \simeq \sqrt{\kappa/m}$ (κ is the force constant that measures the strength of interatomic bonds, and m is the mass involved in the oscillation) as proof that the typical frequencies of lattice vibrations are very much less than the vibration frequencies of simple molecules. It follows, from $\Delta E = \hbar \omega$, that the quantum energy levels of lattice vibrations are far more closely spaced than the vibrational energy levels of gaseous molecules. Thus, it is likely (and is, indeed, the case) that lattice vibrations are not frozen out at room temperature, but, instead, make their full classical contribution to the molar specific heat of the solid. (See Section 5.5.8.)

If the lattice vibrations behave classically then, according to the equipartition theorem (see Section 5.5.5), each normal mode of oscillation has an associated mean energy $k_B T$, in equilibrium at temperature T [($1/2$) $k_B T$ resides in the kinetic energy of the oscillation, and ($1/2$) $k_B T$ resides in the potential energy]. Thus, the internal energy of the solid is

$$U = 3N k_B T = 3\nu RT, \quad (5.488)$$

where $N = \nu N_A$. It follows that the molar heat capacity at constant volume is

$$c_V = \frac{1}{\nu} \left(\frac{\partial U}{\partial T} \right)_V = 3R \quad (5.489)$$

Solid	c_p	Solid	c_p
Copper	24.5	Aluminium	24.4
Silver	25.5	Tin (white)	26.4
Lead	26.4	Sulphur (rhombic)	22.4
Zinc	25.4	Carbon (diamond)	6.1

Table 5.3: Values of c_p (joules/mole/degree) for some solids at $T = 298^\circ \text{K}$.

for solids. This gives a value of 24.9 joules/mole/degree. In fact, at room temperature, most solids (in particular, metals) have heat capacities that lie remarkably close to this value. This fact was discovered experimentally by Pierre Dulong and Alexis Petite at the beginning of the nineteenth century, and was used to make some of the first crude estimates of the molecular weights of solids. (If we know the molar heat capacity of a substance then we can easily work out how much of it corresponds to one mole, and by weighing this amount, and then dividing the result by Avogadro's number, we can then obtain an estimate of the molecular weight.)

Table 5.3 lists the experimental molar heat capacities, c_p , at constant pressure for various solids. The heat capacity at constant volume is somewhat less than the constant pressure value, but not by much, because solids are fairly incompressible. It can be seen that *Dulong and Petite's law* (i.e., that all solids have a molar heat capacities close to 24.9 joules/mole/degree) holds fairly well for metals. However, the law fails badly for diamond. This is not surprising. As is well known, diamond is an extremely hard substance, so its interatomic bonds must be very strong, suggesting that the force constant, κ , is large. Diamond is also a fairly low-density substance, so the mass, m , involved in lattice vibrations is comparatively small. Both these facts suggest that the typical lattice vibration frequency of diamond ($\omega \simeq \sqrt{\kappa/m}$) is high. In fact, the spacing between the different vibrational energy levels (which scales like $\hbar\omega$) is sufficiently large in diamond for the vibrational degrees of freedom to be largely frozen out at room temperature. This accounts for the anomalously low heat capacity of diamond in Table 5.3.

Dulong and Petite's law is essentially a high-temperature limit. We can make a crude model of the behavior of c_V at low temperatures by assuming that all of the normal modes oscillate at the same frequency, ω (say). This approximation was first employed by Einstein in a paper published in 1907. According to Equation (5.487), the solid acts like a set of $3N$ independent oscillators which, making use of Einstein's approximation, all vibrate at the same frequency. We can use the quantum mechanical result (5.390) for the mean energy of a single oscillator to write the internal energy of the solid in the form

$$U = 3N \left[\frac{1}{2} + \frac{1}{\exp(\hbar\omega/k_B T) - 1} \right] \hbar\omega, \quad (5.490)$$

giving

$$c_V \frac{1}{V} \left(\frac{\partial U}{\partial T} \right)_V = 3R \left(\frac{\theta_E}{T} \right)^2 \frac{\exp(\theta_E/T)}{[\exp(\theta_E/T) - 1]^2}. \quad (5.491)$$

Here,

$$\theta_E = \frac{\hbar \omega}{k_B} \quad (5.492)$$

is termed the *Einstein temperature*. If the temperature is sufficiently high that $T \gg \theta_E$ then the previous expression reduces to $c_V = 3R$, after expansion of the exponential functions. Thus, the law of Dulong and Petite is recovered for temperatures significantly in excess of the Einstein temperature. On the other hand, if the temperature is sufficiently low that $T \ll \theta_E$ then the exponential factors appearing in Equation (5.491) become very much larger than unity, giving

$$c_V \simeq 3R \left(\frac{\theta_E}{T} \right)^2 \exp \left(-\frac{\theta_E}{T} \right). \quad (5.493)$$

So, in this simple model, the specific heat approaches zero exponentially as $T \rightarrow 0$.

In reality, the specific heats of solids do not approach zero quite as quickly as suggested by Einstein's model when $T \rightarrow 0$. The experimentally observed low-temperature behavior is more like $c_V \propto T^3$. The reason for this discrepancy is the crude approximation that all normal modes have the same frequency. In fact, long-wavelength modes have lower frequencies than short-wavelength modes, so the former are much harder to freeze out than the latter (because the spacing between quantum energy levels, $\hbar \omega$, is smaller in the former case). The molar heat capacity does not decrease with temperature as rapidly as suggested by Einstein's model because these long-wavelength modes are able to make a significant contribution to the heat capacity, even at very low temperatures. A more realistic model of lattice vibrations was developed by Peter Debye in 1912. In the Debye model, the frequencies of the normal modes of vibration are estimated by treating the solid as an isotropic continuous medium. This approach is reasonable because the only modes that really matter at low temperatures are the long-wavelength modes; more explicitly, those whose wavelengths greatly exceed the interatomic spacing. It is plausible that these modes are not particularly sensitive to the discrete nature of the solid. In other words, they are not sensitive to the fact that the solid is made up of atoms, rather than being continuous.

According to Equation (5.461), the density of sound wave states in a continuous solid is

$$\rho_c(\omega) = \frac{3V}{2\pi^2} \frac{\omega^2}{v_s^3}, \quad (5.494)$$

where v_s is the average sound speed. The Debye approach consists in approximating the actual density of sound wave states, $\rho(\omega)$, by the density in a continuous medium, $\rho_c(\omega)$, not only at low frequencies (long wavelengths) where these should be nearly the same, but also at high frequencies where they may differ substantially. Suppose that we are dealing with a solid consisting of N atoms. We know that there are only $3N$ independent normal modes. It follows that we must cut off the density of states above some critical frequency, ω_D (say), otherwise we will have too many modes. Thus, in the Debye approximation, the density of normal modes takes the form

$$\rho_D(\omega) = \begin{cases} \rho_c(\omega) & \omega \leq \omega_D \\ 0 & \omega > \omega_D \end{cases}. \quad (5.495)$$

Here, ω_D is termed the *Debye frequency*, and is chosen such that the total number of normal modes is $3N$:

$$\int_0^\infty \rho_D(\omega) d\omega = \int_0^{\omega_D} \rho_c(\omega) d\omega = 3N. \quad (5.496)$$

Substituting Equation (5.494) into the previous formula yields

$$\frac{3V}{2\pi^2 v_s^3} \int_0^{\omega_D} \omega^2 d\omega = \frac{V}{2\pi^2 v_s^3} \omega_D^3 = 3N. \quad (5.497)$$

This implies that

$$\omega_D = v_s \left(6\pi^2 \frac{N}{V} \right)^{1/3}. \quad (5.498)$$

Thus, the Debye frequency depends only on the sound speed in the solid, and the number of atoms per unit volume. The wavelength corresponding to the Debye frequency is $2\pi v_s/\omega_D$, which is clearly on the order of the interatomic spacing, $d \simeq (V/N)^{1/3}$. It follows that the cut-off of normal modes whose frequencies exceed the Debye frequency is equivalent to a cut-off of normal modes whose wavelengths are less than the interatomic spacing. Of course, it makes physical sense that such modes should be absent.

We can use the quantum-mechanical expression for the mean energy of a single oscillator, Equation (5.390), to calculate the internal energy associated with lattice vibrations in the Debye approximation. We obtain

$$U = \int_0^\infty \rho_D(\omega) \left[\frac{1}{2} + \frac{1}{\exp(\hbar\omega/k_B T) - 1} \right] \hbar\omega d\omega. \quad (5.499)$$

Hence, the molar heat capacity takes the form

$$c_V = \frac{1}{v} \left(\frac{\partial U}{\partial T} \right)_V = \frac{1}{v k_B T^2} \int_0^\infty \rho_D(\omega) \left\{ \frac{\exp(\hbar\omega/k_B T) \hbar\omega}{[\exp(\hbar\omega/k_B T) - 1]^2} \right\} \hbar\omega d\omega. \quad (5.500)$$

Making use of Equations (5.494) and (5.495), we find that

$$c_V = \frac{k_B}{v} \int_0^{\omega_D} \frac{\exp(\hbar\omega/k_B T) (\hbar\omega/k_B T)^2}{[\exp(\hbar\omega/k_B T) - 1]^2} \frac{3V}{2\pi^2 v_s^3} \omega^2 d\omega, \quad (5.501)$$

giving

$$c_V = \frac{3V k_B}{2\pi^2 v (v_s \hbar/k_B T)^3} \int_0^{\hbar\omega_D/k_B T} \frac{x^4 e^x}{(e^x - 1)^2} dx, \quad (5.502)$$

in terms of the dimensionless variable $x = \hbar\omega/k_B T$. According to Equation (5.498), the volume can be written

$$V = 6\pi^2 N \left(\frac{v_s}{\omega_D} \right)^3, \quad (5.503)$$

so the heat capacity reduces to

$$c_V = 3R f_D \left(\frac{\hbar\omega_D}{k_B T} \right) = 3R f_D \left(\frac{\theta_D}{T} \right), \quad (5.504)$$

Solid	θ_D (low temperature)	θ_D (sound speed)
NaCl	308	320
KCl	230	246
Ag	225	216
Zn	308	305

Table 5.4: Comparison of Debye temperatures (in degrees kelvin) obtained from the low temperature behavior of the heat capacity with those calculated from the sound speed.

where the *Debye function* is defined

$$f_D(y) \equiv \frac{3}{y^3} \int_0^y \frac{x^4 e^x}{(e^x - 1)^2} dx. \quad (5.505)$$

We have also defined the *Debye temperature*, θ_D , as

$$k_B \theta_D = \hbar \omega_D. \quad (5.506)$$

Consider the asymptotic limit in which $T \gg \theta_D$. For small y , we can approximate e^x as $1 + x$ in the integrand of Equation (5.505), so that

$$f_D(y) \rightarrow \frac{3}{y^3} \int_0^y x^2 dx = 1. \quad (5.507)$$

Thus, if the temperature greatly exceeds the Debye temperature then we recover the law of Dulong and Petite that $c_V = 3R$. Consider, now, the asymptotic limit in which $T \ll \theta_D$. For large y ,

$$\int_0^y \frac{x^4 e^x}{(e^x - 1)^2} dx \simeq \int_0^\infty \frac{x^4 e^x}{(e^x - 1)^2} dx = \frac{4\pi^4}{15}. \quad (5.508)$$

The latter integral can be looked up in standard reference books on integrals. Thus, in the low-temperature limit,

$$f_D(y) \rightarrow \frac{4\pi^4}{5} \frac{1}{y^3}, \quad (5.509)$$

which yields

$$c_V \simeq \frac{12\pi^4}{5} R \left(\frac{T}{\theta_D} \right)^3 \quad (5.510)$$

in the limit $T \ll \theta_D$. Note that c_V varies with temperature as T^3 , in accordance with experimental observation.

The fact that c_V goes like T^3 at low temperatures is quite well verified experimentally, although it is sometimes necessary to go to temperatures as low as $0.02 \theta_D$ to obtain this asymptotic behavior. Theoretically, θ_D should be calculable from Equation (5.498) in terms of the sound speed in the solid, and the molar volume. Table 5.4 shows a comparison of Debye temperatures evaluated by this means with temperatures obtained empirically by fitting the law (5.510) to the low-temperature

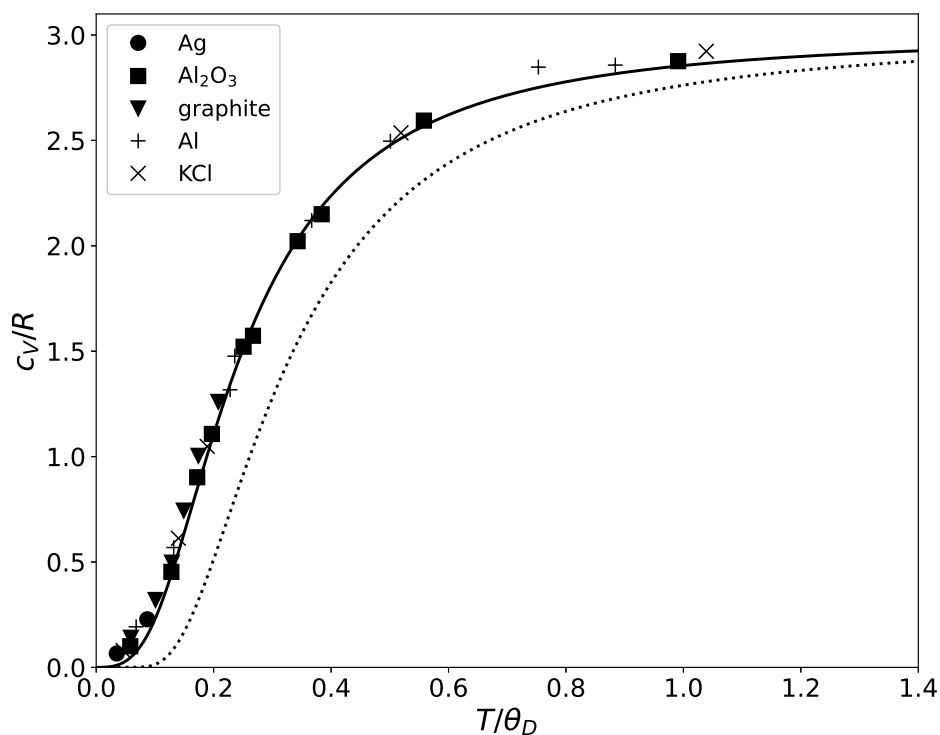


Figure 5.9: The molar heat capacity of various solids. The solid curve shows the prediction of Debye theory. The dotted curve shows the prediction of Einstein theory (assuming that $\theta_E = \theta_D$).

variation of the heat capacity. It can be seen that there is fairly good agreement between the theoretical and empirical Debye temperatures. This suggests that the Debye theory affords a good, though not perfect, representation of the behavior of c_V in solids over the entire temperature range.

Finally, Figure 5.9 shows the actual temperature variation of the molar heat capacities of various solids, as well as that predicted by Debye's theory. The prediction of Einstein's theory is also shown, for the sake of comparison.

5.6.6 Conduction Electrons in Metal

The conduction electrons in a metal are non-localized (i.e., they are not tied to any particular atoms). In conventional metals, each atom contributes a fixed number of such electrons (corresponding to its valency). To a first approximation, it is possible to neglect the mutual interaction of the conduction electrons, because this interaction is largely shielded out by the stationary ions. The conduction electrons can, therefore, be treated as an ideal gas. However, the number density of such electrons in a metal far exceeds the number density of molecules in a conventional gas.

Electrons are subject to the Pauli exclusion principle, according to which a given electron state

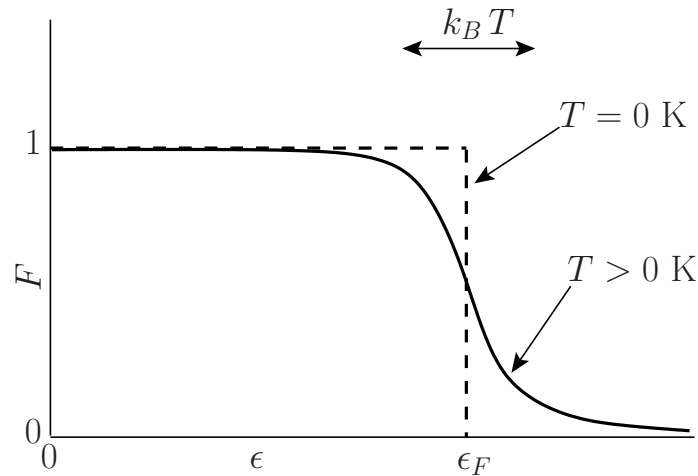


Figure 5.10: The Fermi function.

can either be unoccupied, or singly occupied. (See Section 4.4.3.) The so-called *Fermi energy* is the energy at which electrons are available without doing work. Thus, an electron state of energy ϵ has an available free energy 0 when it is unoccupied, and an available free energy $\epsilon - \epsilon_F$ when it is occupied. According to the Boltzmann distribution (see Section 5.4.7), the relative probabilities of unoccupied and occupied states are thus $P(0) = 1$ and $P(1) = \exp[-(\epsilon - \epsilon_F)/(k_B T)]$, respectively, where T is the temperature. (See Section 5.4.7.) Thus, the mean occupancy number of the state is

$$F(\epsilon) = \frac{0 \times P(0) + 1 \times P(1)}{P(0) + P(1)}, \quad (5.511)$$

which reduces to

$$F(\epsilon) = \frac{1}{\exp[(\epsilon - \epsilon_F)/(k_B T)] + 1}. \quad (5.512)$$

Here, $F(\epsilon)$ is termed the *Fermi function*.

Let us investigate the behavior of the Fermi function as ϵ varies. Here, the energy is measured from its lowest possible value $\epsilon = 0$. The Fermi energy for conduction electrons in a metal is such that $\epsilon_F \gg k_B T$. In this limit, if $0 < \epsilon \ll \epsilon_F$ then $(\epsilon - \epsilon_F)/(k_B T) \ll -1$, so that $F(\epsilon) \simeq 1$. On the other hand, if $\epsilon \gg \epsilon_F$ then $(\epsilon - \epsilon_F)/(k_B T) \gg 1$, so that $F(\epsilon) \simeq \exp[-(\epsilon - \mu)/(k_B T)]$ falls off exponentially with increasing ϵ . Note that $F = 1/2$ when $\epsilon = \epsilon_F$. The transition region in which F goes from a value close to unity to a value close to zero corresponds to an energy interval of order $k_B T$, centered on $\epsilon = \epsilon_F$. In fact, $F = 3/4$ when $\epsilon = \epsilon_F - (\ln 3)k_B T$, and $F = 1/4$ when $\epsilon = \epsilon_F + (\ln 3)k_B T$. The behavior of the Fermi function is illustrated in Figure 5.10.

In the limit as $T \rightarrow 0$, the transition region becomes infinitesimally narrow. In this case, $F = 1$ for $\epsilon \leq \epsilon_F$, and $F = 0$ for $\epsilon > \epsilon_F$, as illustrated in Figure 5.10. This is an obvious result, because when $T = 0$ the conduction electrons attain their lowest energy, or ground-state, configuration. Because the Pauli exclusion principle requires that there be no more than one electron per single-particle quantum state, the lowest energy configuration is obtained by piling electrons into the lowest available unoccupied states, until all of the electrons are used up. Thus, the last electron

added to the pile has a quite considerable energy, $\epsilon = \epsilon_F$, because all of the lower energy states are already occupied. Clearly, the exclusion principle implies that free electrons in a metal possess a large mean energy, even at a temperature of absolute zero.

We can calculate the Fermi energy by equating the number of occupied electron states to the total number of electrons in the metal, N_e . In other words,

$$\int_0^{\infty} F(\epsilon) \rho(\epsilon) d\epsilon = N_e, \quad (5.513)$$

where $\rho(\epsilon)$ is the density of electron states specified in Equation (5.466). Assuming that $\epsilon_F \gg k_B T$, we can make the approximation

$$F(\epsilon) \simeq \begin{cases} 1 & 0 \leq \epsilon \leq \epsilon_F \\ 0 & \epsilon > \epsilon_F \end{cases}. \quad (5.514)$$

Thus, making use of Equation (5.466), we obtain

$$\frac{\sqrt{2} V m_e^{3/2}}{\pi^2 \hbar^3} \int_0^{\epsilon_F} \epsilon^{1/2} d\epsilon = \frac{\sqrt{8} V m_e^{3/2} \epsilon_F^{3/2}}{3 \pi^2 \hbar^3} = N_e, \quad (5.515)$$

which can be rearranged to give

$$\epsilon_F = \frac{\hbar^2}{2 m_e} (3 \pi^2 n_e)^{2/3} \quad (5.516)$$

where $n_e = N_e/V$ is the number density of conduction electrons. The mean electron energy is

$$\langle \epsilon \rangle = \frac{\int_0^{\epsilon_F} \epsilon \epsilon^{1/2} d\epsilon}{\int_0^{\epsilon_F} \epsilon^{1/2} d\epsilon} = \frac{2}{5} \epsilon_F. \quad (5.517)$$

(See Section 4.4.3.)

Copper at room temperature has a number density of conduction electrons of $n_e = 8.4 \times 10^{28} \text{ m}^{-3}$. According to Equation (5.516), the corresponding Fermi energy is

$$\epsilon_F = 7.0 \text{ eV}. \quad (5.518)$$

The associated *Fermi temperature* is

$$\theta_F = \frac{\epsilon_F}{k_B} = 8.1 \times 10^4 \text{ K}. \quad (5.519)$$

Thus, at room temperature, $T = 288 \text{ K}$, we obtain

$$\frac{k_B T}{\epsilon_F} = \frac{\theta_F}{T} \simeq \frac{1}{280}, \quad (5.520)$$

which confirms that $\epsilon_F \gg k_B T$ for conduction electrons in a metal.

Let us crudely approximate the Fermi function at finite temperature in the fashion shown in Figure 5.11. As can be seen from the figure, the proportion of thermally excited electrons is the

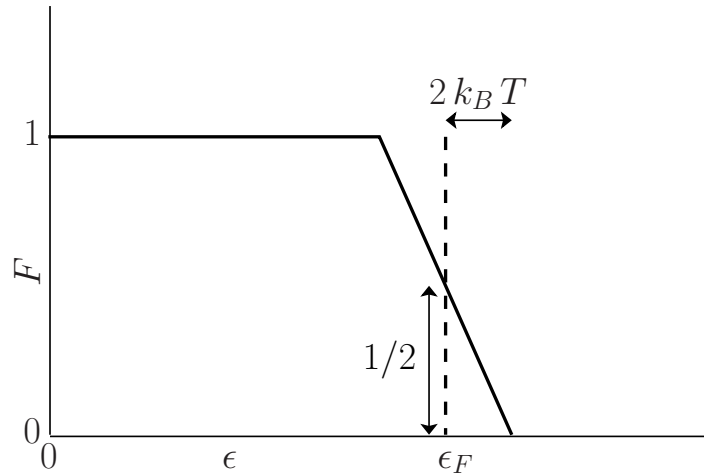


Figure 5.11: Approximate Fermi function.

ratio of the area of a triangle of height $1/2$ and base $2k_B T$ to that of a rectangle of height 1 and base ϵ_F . In other words,

$$\frac{(1/2)(1/2)(2k_B T)}{\epsilon_F} = \frac{k_B T}{2\epsilon_F}. \quad (5.521)$$

Now, the centroid of a right-angled triangle is $1/3$ rd of the distance along its base from the right-angle. Thus, the mean energy of the excited electrons increases by

$$2 \frac{1}{3} 2k_B T = \frac{4}{3} k_B T. \quad (5.522)$$

Hence, the thermal energy per conduction electron is

$$\left(\frac{k_B T}{2\epsilon_F}\right) \left(\frac{4}{3} k_B T\right) = \frac{2k_B^2 T^2}{3\epsilon_F}, \quad (5.523)$$

which implies that the internal energy (i.e., the difference between the energy at a finite temperature and the energy at absolute zero) of the conduction electrons is

$$U \simeq N_e \frac{2k_B^2 T^2}{3\epsilon_F} = \frac{2\nu R k_B T^2}{3\epsilon_F}, \quad (5.524)$$

where ν is the number of moles of electrons. Finally, the molar specific heat of the electrons at constant volume is

$$c_V \simeq \frac{1}{\nu} \left(\frac{\partial U}{\partial T}\right)_V = \frac{4R k_B T}{3\epsilon_F}, \quad (5.525)$$

which can also be written

$$c_V \simeq \left(\frac{3}{2} R\right) \left(\frac{8}{9} \frac{T}{\theta_F}\right). \quad (5.526)$$

The exact result is

$$c_V = \left(\frac{3}{2}R\right) \left(\frac{\pi^2}{3} \frac{T}{\theta_F}\right). \quad (5.527)$$

Thus, we conclude that the contribution of the conduction electrons to the molar specific heat capacity of a metal is proportional to the temperature. However, this contribution is much less than the classical contribution, $(3/2)R$, predicted by the equipartition theorem (see Section 5.5.5), given that each conduction electron possesses three translational degrees of freedom. This is the case because the conduction electrons in a metal are highly degenerate. (See Section 4.4.3.) In fact, Equations (5.520) and (5.527) imply that the contribution of the conduction electrons to the molar specific heat of copper at room temperature is a factor 85 times smaller than the classical contribution.

Using the superscript e to denote the electronic specific heat due to conduction electrons, the molar specific heat of such electrons can be written

$$c_V^{(e)} = \gamma T, \quad (5.528)$$

where γ is a (positive) constant of proportionality. At room temperature, $c_V^{(e)}$ is completely masked by the much larger specific heat, $c_V^{(L)}$, due to lattice vibrations. However, at very low temperatures, $c_V^{(L)} = A T^3$, where A is a (positive) constant of proportionality. (See Section 5.6.5.) Clearly, at low temperatures, $c_V^{(L)} = A T^3$ approaches zero far more rapidly than the electronic specific heat, as T is reduced. Hence, it should be possible to measure the electronic contribution to the molar specific heat at low temperatures.

The total molar specific heat of a metal at low temperatures takes the form

$$c_V = c_V^{(e)} + c_V^{(L)} = \gamma T + A T^3. \quad (5.529)$$

Hence,

$$\frac{c_V}{T} = \gamma + A T^2. \quad (5.530)$$

It follows that a plot of c_V/T versus T^2 should yield a straight-line whose intercept on the vertical axis gives the coefficient γ . Figure 5.12 shows such a plot. The fact that a good straight-line, with a non-zero intercept, is obtained verifies that the temperature dependence of the heat capacity predicted by Equation (5.529) is indeed correct.

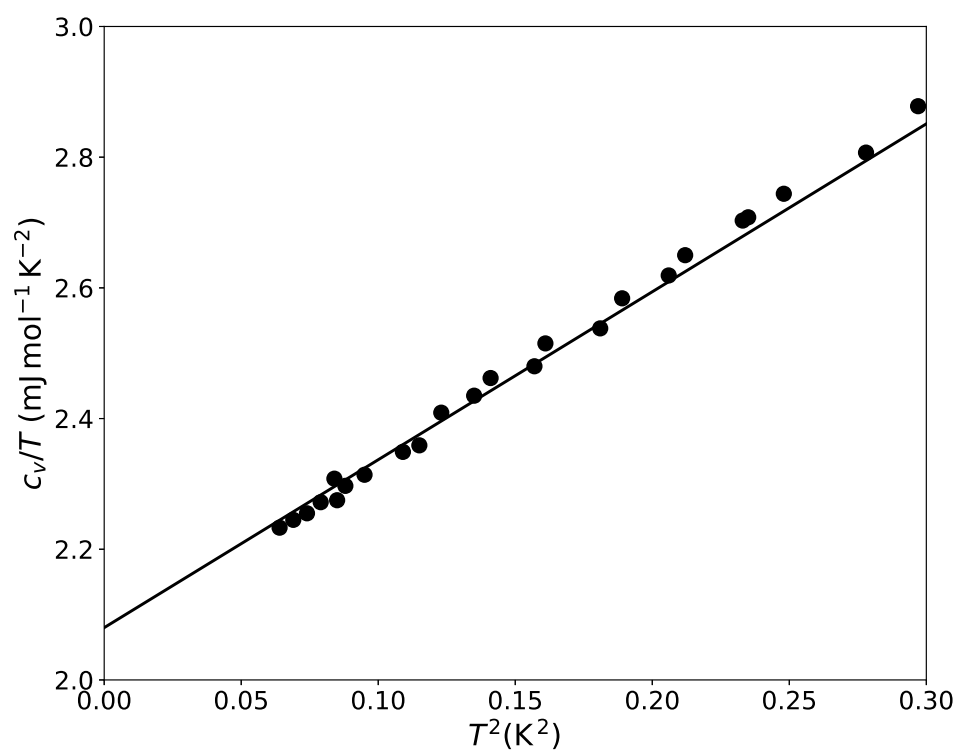


Figure 5.12: The low-temperature heat capacity of potassium, plotted as c_V/T versus T^2 . The straight-line shows the fit $c_V/T = 2.08 + 2.57 T^2$.

Appendix A

Vector Algebra and Vector Calculus

A.1 Introduction

This appendix contains a brief outline of vector algebra and vector calculus. The essential purpose of *vector algebra* is to convert the propositions of Euclidean geometry in three-dimensional space into a convenient algebraic form. *Vector calculus* allows us to define the instantaneous velocity and acceleration of a moving object in three-dimensional space, as well as the work done when such an object travels along a general curved trajectory in a force field. Vector calculus also introduces the concept of a *scalar field*—for example, the potential energy associated with a conservative force field—and a *vector field*—for example, an electric field.

A.2 Scalars and Vectors

Many physical quantities (e.g., mass, energy) are entirely defined by a numerical magnitude (in appropriate units). Such quantities, which have no directional element, are known as *scalars*. Moreover, because scalars can be represented by real numbers, it follows that they obey the familiar laws of ordinary algebra. However, there exists a second class of physical quantities (e.g., velocity, acceleration, force) that are only completely defined when both a numerical magnitude and a direction in space is specified. Such quantities are known as *vectors*. By definition, a vector obeys the same algebra as a displacement in space, and may thus be represented geometrically by

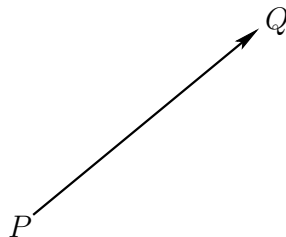


Figure A.1: A vector.

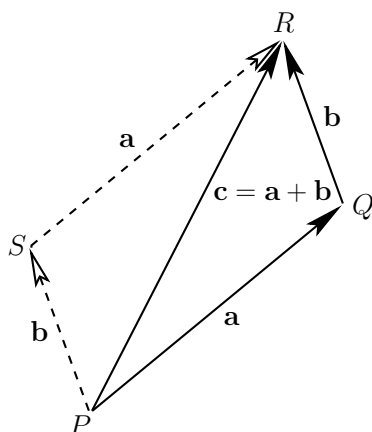


Figure A.2: Vector addition.

a straight-line, \vec{PQ} (say), where the arrow indicates the direction of the displacement (i.e., from point P to point Q). See Figure A.1. The magnitude of the vector is represented by the length of the straight-line.

It is conventional to denote vectors by bold-faced symbols (e.g., \mathbf{a} , \mathbf{F}) and scalars by non-bold-faced symbols (e.g., r , S). The magnitude of a general vector, \mathbf{a} , is denoted $|\mathbf{a}|$, or just a , and is, by definition, always greater than or equal to zero. It is convenient to define a vector with zero magnitude; this is denoted $\mathbf{0}$, and has no direction. Finally, two vectors, \mathbf{a} and \mathbf{b} , are said to be equal when their magnitudes and directions are identical.

A.3 Vector Algebra

Suppose that the displacements \vec{PQ} and \vec{QR} represent the vectors \mathbf{a} and \mathbf{b} , respectively. See Figure A.2. It can be seen that the result of combining these two displacements is to give the net displacement \vec{PR} . Hence, if \vec{PR} represents the vector \mathbf{c} then we can write

$$\mathbf{c} = \mathbf{a} + \mathbf{b}. \quad (\text{A.1})$$

This defines *vector addition*. By completing the parallelogram $PQRS$, we can also see that

$$\vec{PR} = \vec{PQ} + \vec{QR} = \vec{PS} + \vec{SR}. \quad (\text{A.2})$$

However, \vec{PS} has the same length and direction as \vec{QR} , and, thus, represents the same vector, \mathbf{b} . Likewise, \vec{PQ} and \vec{SR} both represent the vector \mathbf{a} . Thus, the previous equation is equivalent to

$$\mathbf{c} = \mathbf{a} + \mathbf{b} = \mathbf{b} + \mathbf{a}. \quad (\text{A.3})$$

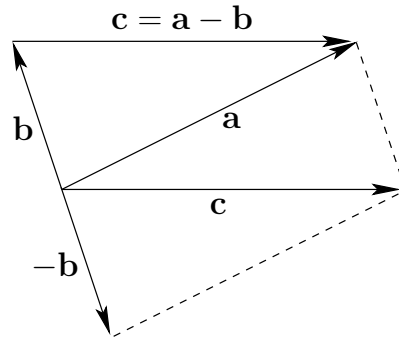


Figure A.3: Vector subtraction.

We conclude that the addition of vectors is *commutative*. It can also be shown that the *associative* law holds; that is,

$$\mathbf{a} + (\mathbf{b} + \mathbf{c}) = (\mathbf{a} + \mathbf{b}) + \mathbf{c}. \quad (\text{A.4})$$

The null vector, $\mathbf{0}$, is represented by a displacement of zero length and arbitrary direction. Because the result of combining such a displacement with a finite length displacement is the same as the latter displacement by itself, it follows that

$$\mathbf{a} + \mathbf{0} = \mathbf{a}, \quad (\text{A.5})$$

where \mathbf{a} is a general vector. The negative of \mathbf{a} is defined as that vector that has the same magnitude, but acts in the opposite direction, and is denoted $-\mathbf{a}$. The sum of \mathbf{a} and $-\mathbf{a}$ is thus the null vector; in other words,

$$\mathbf{a} + (-\mathbf{a}) = \mathbf{0}. \quad (\text{A.6})$$

We can also define the difference of two vectors, \mathbf{a} and \mathbf{b} , as

$$\mathbf{c} = \mathbf{a} - \mathbf{b} = \mathbf{a} + (-\mathbf{b}). \quad (\text{A.7})$$

This definition of *vector subtraction* is illustrated in Figure A.3.

If $n > 0$ is a scalar then the expression $n\mathbf{a}$ denotes a vector whose direction is the same as \mathbf{a} , and whose magnitude is n times that of \mathbf{a} . (This definition becomes obvious when n is an integer.) If n is negative then, because $n\mathbf{a} = |n|(-\mathbf{a})$, it follows that $n\mathbf{a}$ is a vector whose magnitude is $|n|$ times that of \mathbf{a} , and whose direction is opposite to \mathbf{a} . These definitions imply that if n and m are two scalars then

$$n(m\mathbf{a}) = nm\mathbf{a} = m(n\mathbf{a}), \quad (\text{A.8})$$

$$(n+m)\mathbf{a} = n\mathbf{a} + m\mathbf{a}, \quad (\text{A.9})$$

$$n(\mathbf{a} + \mathbf{b}) = n\mathbf{a} + n\mathbf{b}. \quad (\text{A.10})$$

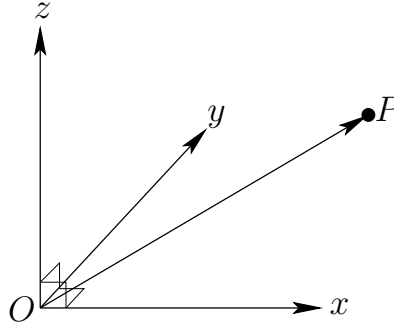


Figure A.4: A right-handed Cartesian coordinate system.

A.4 Cartesian Components of a Vector

Consider a Cartesian coordinate system $Oxyz$ consisting of an origin, O , and three mutually perpendicular coordinate axes, Ox , Oy , and Oz . See Figure A.4. Such a system is said to be *right-handed* if, when looking along the Oz direction, a 90° clockwise rotation about Oz is required to take Ox into Oy . Otherwise, it is said to be left-handed. In physics, it is conventional to always use right-handed coordinate systems.

It is convenient to define unit vectors, \mathbf{e}_x , \mathbf{e}_y , and \mathbf{e}_z , parallel to Ox , Oy , and Oz , respectively. Incidentally, a unit vector is a vector whose magnitude is unity. The position vector, \mathbf{r} , of some general point P whose Cartesian coordinates are (x, y, z) is then given by

$$\mathbf{r} = x\mathbf{e}_x + y\mathbf{e}_y + z\mathbf{e}_z. \quad (\text{A.11})$$

In other words, we can get from O to P by moving a distance x parallel to Ox , then a distance y parallel to Oy , and then a distance z parallel to Oz . Similarly, if \mathbf{a} is an arbitrary vector then

$$\mathbf{a} = a_x\mathbf{e}_x + a_y\mathbf{e}_y + a_z\mathbf{e}_z, \quad (\text{A.12})$$

where a_x , a_y , and a_z are termed the *Cartesian components* of \mathbf{a} . It is conventional to write $\mathbf{a} \equiv (a_x, a_y, a_z)$. It follows that $\mathbf{e}_x \equiv (1, 0, 0)$, $\mathbf{e}_y \equiv (0, 1, 0)$, and $\mathbf{e}_z \equiv (0, 0, 1)$. Of course, $\mathbf{0} \equiv (0, 0, 0)$.

According to the three-dimensional generalization of the Pythagorean theorem, the distance $OP \equiv |\mathbf{r}| = r$ is given by

$$r = \sqrt{x^2 + y^2 + z^2}. \quad (\text{A.13})$$

By analogy, the magnitude of a general vector \mathbf{a} takes the form

$$a = \sqrt{a_x^2 + a_y^2 + a_z^2}. \quad (\text{A.14})$$

If $\mathbf{a} \equiv (a_x, a_y, a_z)$ and $\mathbf{b} \equiv (b_x, b_y, b_z)$ then it is easily demonstrated that

$$\mathbf{a} + \mathbf{b} \equiv (a_x + b_x, a_y + b_y, a_z + b_z). \quad (\text{A.15})$$

Furthermore, if n is a scalar then it is apparent that

$$n\mathbf{a} \equiv (na_x, na_y, na_z). \quad (\text{A.16})$$

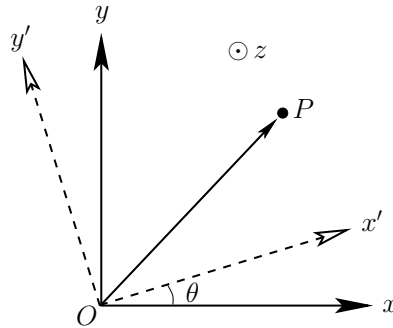


Figure A.5: Rotation of the coordinate axes about Oz .

A.5 Coordinate Transformations

A Cartesian coordinate system allows position and direction in space to be represented in a very convenient manner. Unfortunately, such a coordinate system also introduces arbitrary elements into our analysis. After all, two independent observers might well choose coordinate systems with different origins, and different orientations of the coordinate axes. In general, a given vector \mathbf{a} will have different sets of components in these two coordinate systems. However, the direction and magnitude of \mathbf{a} are the same in both cases. Hence, the two sets of components must be related to one another in a very particular fashion. Actually, because vectors are represented by moveable line elements in space (i.e., in Figure A.2, \vec{PQ} and \vec{SR} represent the same vector), it follows that the components of a general vector are not affected by a simple shift in the origin of a Cartesian coordinate system. On the other hand, the components are modified when the coordinate axes are rotated.

Suppose that we transform to a new coordinate system, $Ox'y'z'$, that has the same origin as $Oxyz$, and is obtained by rotating the coordinate axes of $Oxyz$ through an angle θ about Oz . See Figure A.5. Let the coordinates of a general point P be (x, y, z) in $Oxyz$ and (x', y', z') in $Ox'y'z'$. According to simple trigonometry, these two sets of coordinates are related to one another via the transformation:

$$x' = x \cos \theta + y \sin \theta, \tag{A.17}$$

$$y' = -x \sin \theta + y \cos \theta, \tag{A.18}$$

$$z' = z. \tag{A.19}$$

Consider the vector displacement $\mathbf{r} \equiv \vec{OP}$. Note that this displacement is represented by the same symbol, \mathbf{r} , in both coordinate systems, because the magnitude and direction of \mathbf{r} are manifestly independent of the orientation of the coordinate axes. The coordinates of \mathbf{r} do depend on the orientation of the axes; that is, $\mathbf{r} \equiv (x, y, z)$ in $Oxyz$, and $\mathbf{r} \equiv (x', y', z')$ in $Ox'y'z'$. However, they must depend in a very specific manner [i.e., Equations (A.17)–(A.19)] that preserves the magnitude and direction of \mathbf{r} .

The components of a general vector \mathbf{a} transform in an analogous manner to Equations (A.17)–(A.19); that is,

$$a_{x'} = a_x \cos \theta + a_y \sin \theta, \quad (\text{A.20})$$

$$a_{y'} = -a_x \sin \theta + a_y \cos \theta, \quad (\text{A.21})$$

$$a_{z'} = a_z. \quad (\text{A.22})$$

Moreover, there are similar transformation rules for rotation about Ox and Oy . Equations (A.20)–(A.22) effectively constitute the definition of a vector; in other words, the three quantities (a_x, a_y, a_z) are the components of a vector provided that they transform under rotation of the coordinate axes about Oz in accordance with Equations (A.20)–(A.22). (And also transform correctly under rotation about Ox and Oy). Conversely, (a_x, a_y, a_z) cannot be the components of a vector if they do not transform in accordance with Equations (A.20)–(A.22). Of course, scalar quantities are invariant under rotation of the coordinate axes. Thus, the individual components of a vector (a_x , say) are real numbers, but they are not scalars. Displacement vectors, and all vectors derived from displacements (e.g., velocity, acceleration), automatically satisfy Equations (A.20)–(A.22). There are, however, other physical quantities that have both magnitude and direction, but that are not obviously related to displacements. We need to check carefully to see whether these quantities are really vectors. (See Section A.9.)

A.6 Scalar Product

A scalar quantity is invariant under all possible rotational transformations. The individual components of a vector are not scalars because they change under transformation. Can we form a scalar out of some combination of the components of one, or more, vectors? Suppose that we were to define the “percent” product,

$$\mathbf{a} \% \mathbf{b} \equiv a_x b_z + a_y b_x + a_z b_y = \text{scalar number}, \quad (\text{A.23})$$

for general vectors \mathbf{a} and \mathbf{b} . Is $\mathbf{a} \% \mathbf{b}$ invariant under transformation, as must be the case if it is a scalar number? Let us consider an example. Suppose that $\mathbf{a} \equiv (0, 1, 0)$ and $\mathbf{b} \equiv (1, 0, 0)$. It is easily seen that $\mathbf{a} \% \mathbf{b} = 1$. Let us now rotate the coordinate axes through 45° about Oz . In the new coordinate system, $\mathbf{a} \equiv (1/\sqrt{2}, 1/\sqrt{2}, 0)$ and $\mathbf{b} \equiv (1/\sqrt{2}, -1/\sqrt{2}, 0)$, giving $\mathbf{a} \% \mathbf{b} = 1/2$. Clearly, $\mathbf{a} \% \mathbf{b}$ is not invariant under rotational transformation, so the previous definition is a bad one.

Consider, now, the *dot product* or *scalar product*:

$$\mathbf{a} \cdot \mathbf{b} \equiv a_x b_x + a_y b_y + a_z b_z = \text{scalar number}. \quad (\text{A.24})$$

Let us rotate the coordinate axes through θ degrees about Oz . According to Equations (A.20)–(A.22), $\mathbf{a} \cdot \mathbf{b}$ takes the form

$$\mathbf{a} \cdot \mathbf{b} = (a_x \cos \theta + a_y \sin \theta)(b_x \cos \theta + b_y \sin \theta)$$

$$\begin{aligned}
&+ (-a_x \sin \theta + a_y \cos \theta)(-b_x \sin \theta + b_y \cos \theta) + a_z b_z \\
&= a_x b_x + a_y b_y + a_z b_z
\end{aligned} \tag{A.25}$$

in the new coordinate system. Thus, $\mathbf{a} \cdot \mathbf{b}$ is invariant under rotation about Oz . It can easily be shown that it is also invariant under rotation about Ox and Oy . We conclude that $\mathbf{a} \cdot \mathbf{b}$ is a true scalar, and that the previous definition is a good one. Incidentally, $\mathbf{a} \cdot \mathbf{b}$ is the only simple combination of the components of two vectors that transforms like a scalar. It is easily shown that the dot product is commutative and distributive; that is,

$$\begin{aligned}
\mathbf{a} \cdot \mathbf{b} &= \mathbf{b} \cdot \mathbf{a}, \\
\mathbf{a} \cdot (\mathbf{b} + \mathbf{c}) &= \mathbf{a} \cdot \mathbf{b} + \mathbf{a} \cdot \mathbf{c}.
\end{aligned} \tag{A.26}$$

The associative property is meaningless for the dot product, because we cannot have $(\mathbf{a} \cdot \mathbf{b}) \cdot \mathbf{c}$, because $\mathbf{a} \cdot \mathbf{b}$ is scalar.

We have shown that the dot product $\mathbf{a} \cdot \mathbf{b}$ is coordinate independent. But what is the geometric significance of this? Well, in the special case where $\mathbf{a} = \mathbf{b}$, we get

$$\mathbf{a} \cdot \mathbf{b} = a_x^2 + a_y^2 + a_z^2 = |\mathbf{a}|^2 = a^2. \tag{A.27}$$

So, the invariance of $\mathbf{a} \cdot \mathbf{a}$ is equivalent to the invariance of the magnitude of vector \mathbf{a} under transformation.

Let us now investigate the general case. The length squared of AB in the vector triangle shown in Figure A.6 is

$$(\mathbf{b} - \mathbf{a}) \cdot (\mathbf{b} - \mathbf{a}) = |\mathbf{a}|^2 + |\mathbf{b}|^2 - 2 \mathbf{a} \cdot \mathbf{b}. \tag{A.28}$$

However, according to the ‘‘cosine rule’’ of trigonometry,

$$(AB)^2 = (OA)^2 + (OB)^2 - 2(OA)(OB) \cos \theta, \tag{A.29}$$

where (AB) denotes the length of side AB . It follows that

$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}| |\mathbf{b}| \cos \theta. \tag{A.30}$$

In this case, the invariance of $\mathbf{a} \cdot \mathbf{b}$ under transformation is equivalent to the invariance of the angle subtended between the two vectors. Note that if $\mathbf{a} \cdot \mathbf{b} = 0$ then either $|\mathbf{a}| = 0$, $|\mathbf{b}| = 0$, or the vectors \mathbf{a} and \mathbf{b} are mutually perpendicular. The angle subtended between two vectors can easily be obtained from the dot product; in fact,

$$\cos \theta = \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{a}| |\mathbf{b}|}. \tag{A.31}$$

The work W performed by a constant force \mathbf{F} that moves an object through a displacement \mathbf{r} is the product of the magnitude of \mathbf{F} and the displacement in the direction of \mathbf{F} . If the angle subtended between \mathbf{F} and \mathbf{r} is θ then

$$W = |\mathbf{F}| (|\mathbf{r}| \cos \theta) = \mathbf{F} \cdot \mathbf{r}. \tag{A.32}$$

The work dW performed by a non-constant force \mathbf{f} that moves an object through an infinitesimal displacement $d\mathbf{r}$ in a time interval dt is $dW = \mathbf{f} \cdot d\mathbf{r}$. Thus, the rate at which the force does work on the object, which is usually referred to as the power, is $P = dW/dt = \mathbf{f} \cdot d\mathbf{r}/dt$, or $P = \mathbf{f} \cdot \mathbf{v}$, where $\mathbf{v} = d\mathbf{r}/dt$ is the object’s instantaneous velocity.

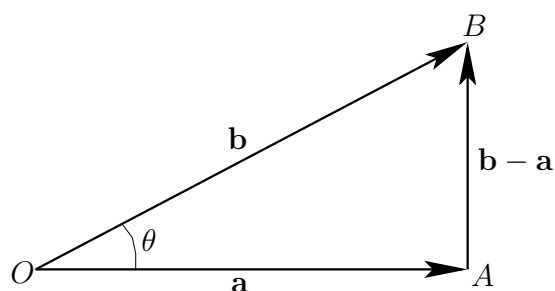


Figure A.6: A vector triangle.

A.7 Vector Area

Suppose that we have planar surface of scalar area S . We can define a vector area \mathbf{S} whose magnitude is S , and whose direction is perpendicular to the plane, in the sense determined by a right-hand circulation rule (see Section A.8) applied to the rim, assuming that a direction of circulation around the rim is specified. (See Figure A.7.) This quantity clearly possesses both magnitude and direction. But is it a true vector? We know that if the normal to the surface makes an angle α_x with the x -axis then the area seen looking along the x -direction is $S \cos \alpha_x$. This is the x -component of \mathbf{S} (because $S_x = \mathbf{e}_x \cdot \mathbf{S} = \mathbf{e}_x \cdot \mathbf{n}S = \cos \alpha_x S$, where \mathbf{n} is the unit normal to the surface). Similarly, if the normal makes an angle α_y with the y -axis then the area seen looking along the y -direction is $S \cos \alpha_y$. This is the y -component of \mathbf{S} . If we limit ourselves to a surface whose normal is perpendicular to the z -direction then $\alpha_x = \pi/2 - \alpha_y = \alpha$. It follows that $\mathbf{S} = S (\cos \alpha, \sin \alpha, 0)$. If we rotate the basis about the z -axis by θ degrees, which is equivalent to rotating the normal to the surface about the z -axis by $-\theta$ degrees, so that $\alpha \rightarrow \alpha - \theta$, then

$$S_{x'} = S \cos(\alpha - \theta) = S \cos \alpha \cos \theta + S \sin \alpha \sin \theta = S_x \cos \theta + S_y \sin \theta, \quad (\text{A.33})$$

which is the correct transformation rule for the x -component of a vector. The other components transform correctly as well. This proves both that a vector area is a true vector, and that the components of a vector area are the projected areas seen looking down the coordinate axes.

According to the vector addition theorem, the projected area of two plane surfaces, joined together at a line, looking along the x -direction (say) is the x -component of the resultant of the vector areas of the two surfaces. Likewise, for many joined-up plane areas, the net area seen looking down the x -axis, which is the same as the area of the outer rim seen looking down the x -axis, is the x -component of the resultant of all the vector areas: that is,

$$\mathbf{S} = \sum_i \mathbf{S}_i. \quad (\text{A.34})$$

If we approach a limit, by letting the number of plane facets increase, and their areas reduce, then we obtain a continuous surface denoted by the resultant vector area

$$\mathbf{S} = \sum_i \delta \mathbf{S}_i. \quad (\text{A.35})$$

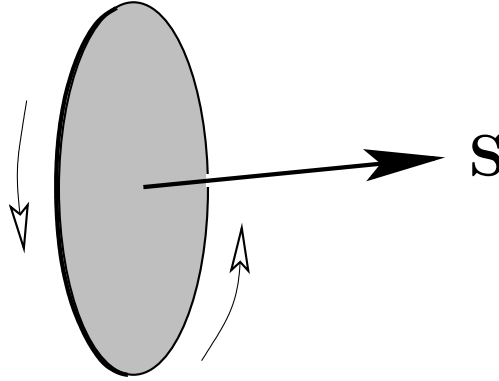


Figure A.7: A vector area.

It is clear that the area of the rim seen looking down the x -axis is just S_x . Similarly, for the areas of the rim seen looking down the other coordinate axes. Note that it is the rim of the surface that determines the vector area, rather than the nature of the surface spanning the rim. So, two different surfaces sharing the same rim both possess the same vector area.

In conclusion, a loop (not all in one plane) has a vector area \mathbf{S} which is the resultant of the component vector areas of any surface ending on the loop. The components of \mathbf{S} are the areas of the loop seen looking down the coordinate axes. As a corollary, a closed surface has $\mathbf{S} = \mathbf{0}$, because it does not possess a rim.

A.8 Vector Product

We have discovered how to construct a scalar from the components of two general vectors \mathbf{a} and \mathbf{b} . Can we also construct a vector that is not just a linear combination of \mathbf{a} and \mathbf{b} ? Consider the following definition:

$$\mathbf{a} * \mathbf{b} \equiv (a_x b_x, a_y b_y, a_z b_z). \tag{A.36}$$

Is $\mathbf{a} * \mathbf{b}$ a proper vector? Suppose that $\mathbf{a} = (0, 1, 0)$, $\mathbf{b} = (1, 0, 0)$. In this case, $\mathbf{a} * \mathbf{b} = \mathbf{0}$. However, if we rotate the coordinate axes through 45° about O_z then $\mathbf{a} = (1/\sqrt{2}, 1/\sqrt{2}, 0)$, $\mathbf{b} = (1/\sqrt{2}, -1/\sqrt{2}, 0)$, and $\mathbf{a} * \mathbf{b} = (1/2, -1/2, 0)$. Thus, $\mathbf{a} * \mathbf{b}$ does not transform like a vector, because its magnitude depends on the choice of axes. So, the previous definition is a bad one.

Consider, now, the *cross product* or *vector product*:

$$\mathbf{a} \times \mathbf{b} \equiv (a_y b_z - a_z b_y, a_z b_x - a_x b_z, a_x b_y - a_y b_x) = \mathbf{c}. \tag{A.37}$$

Does this rather unlikely combination transform like a vector? Let us try rotating the coordinate axes through an angle θ about O_z using Equations (A.20)–(A.22). In the new coordinate system,

$$\begin{aligned} c_{x'} &= (-a_x \sin \theta + a_y \cos \theta) b_z - a_z (-b_x \sin \theta + b_y \cos \theta) \\ &= (a_y b_z - a_z b_y) \cos \theta + (a_z b_x - a_x b_z) \sin \theta \\ &= c_x \cos \theta + c_y \sin \theta. \end{aligned} \tag{A.38}$$

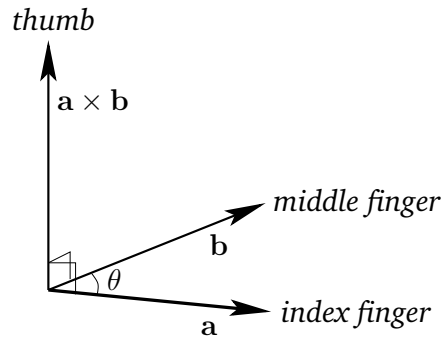


Figure A.8: The right-hand rule for cross products. Here, θ is less than 180° .

Thus, the x -component of $\mathbf{a} \times \mathbf{b}$ transforms correctly. It can easily be shown that the other components transform correctly as well, and that all components also transform correctly under rotation about Ox and Oy . Thus, $\mathbf{a} \times \mathbf{b}$ is a proper vector. Incidentally, $\mathbf{a} \times \mathbf{b}$ is the only simple combination of the components of two vectors that transforms like a vector (and is non-coplanar with \mathbf{a} and \mathbf{b}). The cross product is *anticommutative*,

$$\mathbf{a} \times \mathbf{b} = -\mathbf{b} \times \mathbf{a}, \quad (\text{A.39})$$

distributive,

$$\mathbf{a} \times (\mathbf{b} + \mathbf{c}) = \mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c}, \quad (\text{A.40})$$

but is not associative,

$$\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) \neq (\mathbf{a} \times \mathbf{b}) \times \mathbf{c}. \quad (\text{A.41})$$

The cross product transforms like a vector, which means that it must have a well-defined direction and magnitude. We can show that $\mathbf{a} \times \mathbf{b}$ is perpendicular to both \mathbf{a} and \mathbf{b} . Consider $\mathbf{a} \cdot \mathbf{a} \times \mathbf{b}$. If this is zero then the cross product must be perpendicular to \mathbf{a} . Now,

$$\begin{aligned} \mathbf{a} \cdot \mathbf{a} \times \mathbf{b} &= a_x(a_y b_z - a_z b_y) + a_y(a_z b_x - a_x b_z) + a_z(a_x b_y - a_y b_x) \\ &= 0. \end{aligned} \quad (\text{A.42})$$

Therefore, $\mathbf{a} \times \mathbf{b}$ is perpendicular to \mathbf{a} . Likewise, it can be demonstrated that $\mathbf{a} \times \mathbf{b}$ is perpendicular to \mathbf{b} . The vectors \mathbf{a} , \mathbf{b} , and $\mathbf{a} \times \mathbf{b}$ form a right-handed set, like the unit vectors \mathbf{e}_x , \mathbf{e}_y , and \mathbf{e}_z . In fact, $\mathbf{e}_x \times \mathbf{e}_y = \mathbf{e}_z$. This defines a unique direction for $\mathbf{a} \times \mathbf{b}$, which is obtained from a right-hand rule. See Figure A.8.

Let us now evaluate the magnitude of $\mathbf{a} \times \mathbf{b}$. We have

$$\begin{aligned} (\mathbf{a} \times \mathbf{b})^2 &= (a_y b_z - a_z b_y)^2 + (a_z b_x - a_x b_z)^2 + (a_x b_y - a_y b_x)^2 \\ &= (a_x^2 + a_y^2 + a_z^2)(b_x^2 + b_y^2 + b_z^2) - (a_x b_x + a_y b_y + a_z b_z)^2 \\ &= |\mathbf{a}|^2 |\mathbf{b}|^2 - (\mathbf{a} \cdot \mathbf{b})^2 \\ &= |\mathbf{a}|^2 |\mathbf{b}|^2 - |\mathbf{a}|^2 |\mathbf{b}|^2 \cos^2 \theta = |\mathbf{a}|^2 |\mathbf{b}|^2 \sin^2 \theta. \end{aligned} \quad (\text{A.43})$$

Thus,

$$|\mathbf{a} \times \mathbf{b}| = |\mathbf{a}| |\mathbf{b}| \sin \theta, \quad (\text{A.44})$$

where θ is the angle subtended between \mathbf{a} and \mathbf{b} . Clearly, $\mathbf{a} \times \mathbf{a} = \mathbf{0}$ for any vector, because θ is always zero in this case. Also, if $\mathbf{a} \times \mathbf{b} = \mathbf{0}$ then either $|\mathbf{a}| = 0$, $|\mathbf{b}| = 0$, or \mathbf{b} is parallel (or antiparallel) to \mathbf{a} .

Consider the parallelogram defined by the vectors \mathbf{a} and \mathbf{b} . See Figure A.9. The scalar area of the parallelogram is $ab \sin \theta$. By convention, the *vector area* has the magnitude of the scalar area, and is normal to the plane of the parallelogram, in the sense obtained from a right-hand circulation rule by rotating \mathbf{a} on to \mathbf{b} (through an acute angle); that is, if the fingers of the right-hand circulate in the direction of rotation then the thumb of the right-hand indicates the direction of the vector area. So, the vector area is coming out of the page in Figure A.9. It follows that

$$\mathbf{S} = \mathbf{a} \times \mathbf{b}, \quad (\text{A.45})$$

Suppose that a force \mathbf{F} is applied at position \mathbf{r} . See Figure A.10. The torque about the origin O is the product of the magnitude of the force and the length of the lever arm OQ . Thus, the magnitude of the torque is $|\mathbf{F}| |\mathbf{r}| \sin \theta$. The direction of the torque is conventionally defined as the direction of the axis through O about which the force tries to rotate objects, in the sense determined by a right-hand circulation rule. Hence, the torque is out of the page in Figure A.10. It follows that the vector torque is given by

$$\boldsymbol{\tau} = \mathbf{r} \times \mathbf{F}. \quad (\text{A.46})$$

The angular momentum, \mathbf{l} , of a particle of linear momentum \mathbf{p} and position vector \mathbf{r} is simply defined as the moment of its momentum about the origin: that is,

$$\mathbf{l} = \mathbf{r} \times \mathbf{p}. \quad (\text{A.47})$$

A.9 Rotation

Let us try to define a rotation vector $\boldsymbol{\theta}$ whose magnitude is the angle of the rotation, θ , and whose direction is parallel to the axis of rotation, in the sense determined by a right-hand circulation rule. Unfortunately, this is not a good vector. The problem is that the addition of rotations is not commutative, whereas vector addition is commutative. Figure A.11 shows the effect of applying

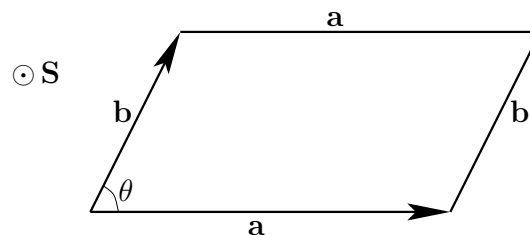


Figure A.9: A vector parallelogram.

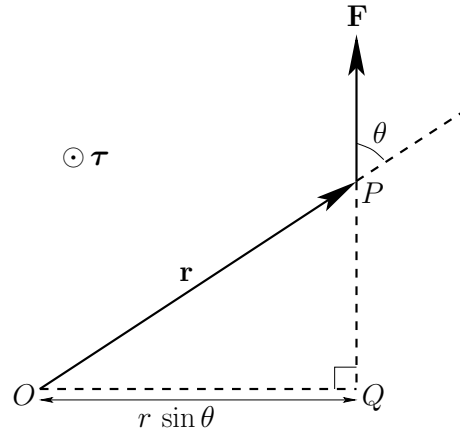


Figure A.10: A torque.

two successive 90° rotations, one about Ox , and the other about the Oz , to a standard six-sided die. In the left-hand case, the z -rotation is applied before the x -rotation, and vice versa in the right-hand case. It can be seen that the die ends up in two completely different states. In other words, the z -rotation plus the x -rotation does not equal the x -rotation plus the z -rotation. This non-commuting algebra cannot be represented by vectors. So, although rotations have a well-defined magnitude and direction, they are not vector quantities.

But, this is not quite the end of the story. Suppose that we take a general vector \mathbf{a} and rotate it about Oz by a small angle $\delta\theta_z$. This is equivalent to rotating the coordinate axes about the Oz by $-\delta\theta_z$. According to Equations (A.20)–(A.22), we have

$$\mathbf{a}' \simeq \mathbf{a} + \delta\theta_z \mathbf{e}_z \times \mathbf{a}, \quad (\text{A.48})$$

where use has been made of the small angle approximations $\sin \theta \simeq \theta$ and $\cos \theta \simeq 1$. The previous equation can easily be generalized to allow small rotations about Ox and Oy by $\delta\theta_x$ and $\delta\theta_y$, respectively. We find that

$$\mathbf{a}' \simeq \mathbf{a} + \delta\boldsymbol{\theta} \times \mathbf{a}, \quad (\text{A.49})$$

where

$$\delta\boldsymbol{\theta} = \delta\theta_x \mathbf{e}_x + \delta\theta_y \mathbf{e}_y + \delta\theta_z \mathbf{e}_z. \quad (\text{A.50})$$

Clearly, we can define a rotation vector, $\delta\boldsymbol{\theta}$, but it only works for small angle rotations (i.e., sufficiently small that the small-angle approximations of sine and cosine are good). According to the previous equation, a small z -rotation plus a small x -rotation is (approximately) equal to the two rotations applied in the opposite order. The fact that infinitesimal rotation is a vector implies that angular velocity,

$$\boldsymbol{\omega} = \lim_{\delta t \rightarrow 0} \frac{\delta\boldsymbol{\theta}}{\delta t}, \quad (\text{A.51})$$

must be a vector as well. Also, if \mathbf{a}' is interpreted as $\mathbf{a}(t + \delta t)$ in Equation (A.49) then it follows that the equation of motion of a vector that precesses about the origin with some angular velocity $\boldsymbol{\omega}$ is

$$\frac{d\mathbf{a}}{dt} = \boldsymbol{\omega} \times \mathbf{a}. \quad (\text{A.52})$$

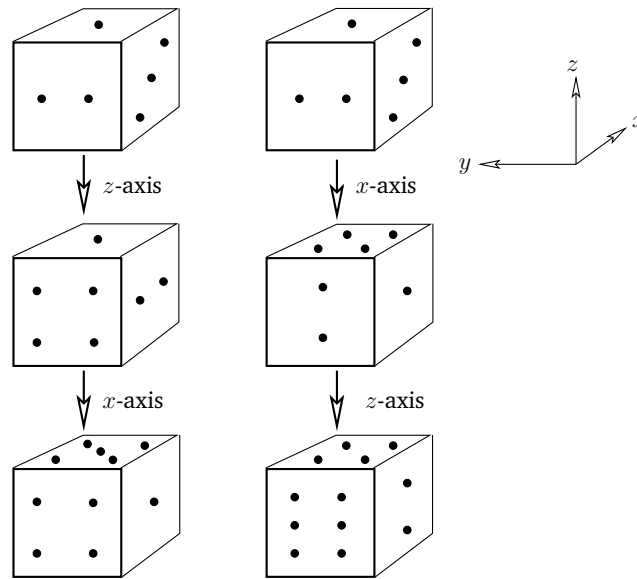


Figure A.11: Effect of successive rotations about perpendicular axes on a six-sided die.

A.10 Scalar Triple Product

Consider three vectors \mathbf{a} , \mathbf{b} , and \mathbf{c} . The *scalar triple product* is defined $\mathbf{a} \cdot \mathbf{b} \times \mathbf{c}$. Now, $\mathbf{b} \times \mathbf{c}$ is the vector area of the parallelogram defined by \mathbf{b} and \mathbf{c} . So, $\mathbf{a} \cdot \mathbf{b} \times \mathbf{c}$ is the scalar area of this parallelogram multiplied by the component of \mathbf{a} in the direction of its normal. It follows that $\mathbf{a} \cdot \mathbf{b} \times \mathbf{c}$ is the volume of the parallelepiped defined by vectors \mathbf{a} , \mathbf{b} , and \mathbf{c} . See Figure A.12. This volume is independent of how the triple product is formed from \mathbf{a} , \mathbf{b} , and \mathbf{c} , except that

$$\mathbf{a} \cdot \mathbf{b} \times \mathbf{c} = -\mathbf{a} \cdot \mathbf{c} \times \mathbf{b}. \quad (\text{A.53})$$

So, the “volume” is positive if \mathbf{a} , \mathbf{b} , and \mathbf{c} form a right-handed set (i.e., if \mathbf{a} lies above the plane of \mathbf{b} and \mathbf{c} , in the sense determined from a right-hand circulation rule by rotating \mathbf{b} onto \mathbf{c}) and negative if they form a left-handed set. The triple product is unchanged if the dot and cross product operators are interchanged,

$$\mathbf{a} \cdot \mathbf{b} \times \mathbf{c} = \mathbf{a} \times \mathbf{b} \cdot \mathbf{c}. \quad (\text{A.54})$$

The triple product is also invariant under any cyclic permutation of \mathbf{a} , \mathbf{b} , and \mathbf{c} ,

$$\mathbf{a} \cdot \mathbf{b} \times \mathbf{c} = \mathbf{b} \cdot \mathbf{c} \times \mathbf{a} = \mathbf{c} \cdot \mathbf{a} \times \mathbf{b}, \quad (\text{A.55})$$

but any anti-cyclic permutation causes it to change sign,

$$\mathbf{a} \cdot \mathbf{b} \times \mathbf{c} = -\mathbf{b} \cdot \mathbf{a} \times \mathbf{c}. \quad (\text{A.56})$$

The scalar triple product is zero if any two of \mathbf{a} , \mathbf{b} , and \mathbf{c} are parallel, or if \mathbf{a} , \mathbf{b} , and \mathbf{c} are coplanar.

If \mathbf{a} , \mathbf{b} , and \mathbf{c} are non-coplanar then any vector \mathbf{r} can be written in terms of them; that is,

$$\mathbf{r} = \alpha \mathbf{a} + \beta \mathbf{b} + \gamma \mathbf{c}. \quad (\text{A.57})$$

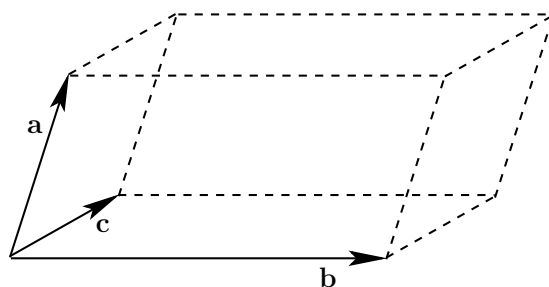


Figure A.12: A vector parallelepiped.

Forming the dot product of this equation with $\mathbf{b} \times \mathbf{c}$, we then obtain

$$\mathbf{r} \cdot \mathbf{b} \times \mathbf{c} = \alpha \mathbf{a} \cdot \mathbf{b} \times \mathbf{c}, \quad (\text{A.58})$$

so

$$\alpha = \frac{\mathbf{r} \cdot \mathbf{b} \times \mathbf{c}}{\mathbf{a} \cdot \mathbf{b} \times \mathbf{c}}. \quad (\text{A.59})$$

Analogous expressions can be written for β and γ . The parameters α , β , and γ are uniquely determined provided $\mathbf{a} \cdot \mathbf{b} \times \mathbf{c} \neq 0$; that is, provided that the three vectors are non-coplanar.

A.11 Vector Triple Product

For three vectors \mathbf{a} , \mathbf{b} , and \mathbf{c} , the *vector triple product* is defined $\mathbf{a} \times (\mathbf{b} \times \mathbf{c})$. The brackets are important because $\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) \neq (\mathbf{a} \times \mathbf{b}) \times \mathbf{c}$. In fact, it can be demonstrated that

$$\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) \equiv (\mathbf{a} \cdot \mathbf{c}) \mathbf{b} - (\mathbf{a} \cdot \mathbf{b}) \mathbf{c} \quad (\text{A.60})$$

and

$$(\mathbf{a} \times \mathbf{b}) \times \mathbf{c} \equiv (\mathbf{a} \cdot \mathbf{c}) \mathbf{b} - (\mathbf{b} \cdot \mathbf{c}) \mathbf{a}. \quad (\text{A.61})$$

Let us try to prove the first of the previous theorems. The left-hand side and the right-hand side are both proper vectors, so if we can prove this result in one particular coordinate system then it must be true in general. Let us take convenient axes such that Ox lies along \mathbf{b} , and \mathbf{c} lies in the x - y plane. It follows that $\mathbf{b} \equiv (b_x, 0, 0)$, $\mathbf{c} \equiv (c_x, c_y, 0)$, and $\mathbf{a} \equiv (a_x, a_y, a_z)$. The vector $\mathbf{b} \times \mathbf{c}$ is directed along Oz ; that is, $\mathbf{b} \times \mathbf{c} \equiv (0, 0, b_x c_y)$. Hence, $\mathbf{a} \times (\mathbf{b} \times \mathbf{c})$ lies in the x - y plane; that is, $\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) \equiv (a_y b_x c_y, -a_x b_x c_y, 0)$. This is the left-hand side of Equation (A.60) in our convenient coordinate system. To evaluate the right-hand side, we need $\mathbf{a} \cdot \mathbf{c} = a_x c_x + a_y c_y$ and $\mathbf{a} \cdot \mathbf{b} = a_x b_x$. It follows that the right-hand side is

$$\begin{aligned} \text{RHS} &= ([a_x c_x + a_y c_y] b_x, 0, 0) - (a_x b_x c_x, a_x b_x c_y, 0) \\ &= (a_y c_y b_x, -a_x b_x c_y, 0) = \text{LHS}, \end{aligned} \quad (\text{A.62})$$

which proves the theorem.

A.12 Vector Calculus

Suppose that vector \mathbf{a} varies with time, so that $\mathbf{a} = \mathbf{a}(t)$. The time derivative of the vector is defined

$$\frac{d\mathbf{a}}{dt} = \lim_{\delta t \rightarrow 0} \left[\frac{\mathbf{a}(t + \delta t) - \mathbf{a}(t)}{\delta t} \right]. \quad (\text{A.63})$$

When written out in component form this becomes

$$\frac{d\mathbf{a}}{dt} \equiv \left(\frac{da_x}{dt}, \frac{da_y}{dt}, \frac{da_z}{dt} \right). \quad (\text{A.64})$$

Suppose that \mathbf{a} is, in fact, the product of a scalar $\phi(t)$ and another vector $\mathbf{b}(t)$. What now is the time derivative of \mathbf{a} ? We have

$$\frac{da_x}{dt} = \frac{d}{dt}(\phi b_x) = \frac{d\phi}{dt} b_x + \phi \frac{db_x}{dt}, \quad (\text{A.65})$$

which implies that

$$\frac{d\mathbf{a}}{dt} = \frac{d\phi}{dt} \mathbf{b} + \phi \frac{d\mathbf{b}}{dt}. \quad (\text{A.66})$$

Moreover, it is easily demonstrated that

$$\frac{d}{dt}(\mathbf{a} \cdot \mathbf{b}) = \frac{d\mathbf{a}}{dt} \cdot \mathbf{b} + \mathbf{a} \cdot \frac{d\mathbf{b}}{dt}, \quad (\text{A.67})$$

and

$$\frac{d}{dt}(\mathbf{a} \times \mathbf{b}) = \frac{d\mathbf{a}}{dt} \times \mathbf{b} + \mathbf{a} \times \frac{d\mathbf{b}}{dt}. \quad (\text{A.68})$$

Hence, it can be seen that the laws of vector differentiation are analogous to those in conventional calculus.

A.13 Line Integrals

Consider a two-dimensional function $f(x, y)$ that is defined for all x and y . What is meant by the integral of f along a given curve joining the points P and Q in the x - y plane? Well, we first draw out f as a function of length l along the path. See Figure A.13. The integral is then simply given by

$$\int_P^Q f(x, y) dl = \text{Area under the curve}, \quad (\text{A.69})$$

where $dl = (dx^2 + dy^2)^{1/2}$.

As an example of this, consider the integral of $f(x, y) = xy^2$ between P and Q along the two routes indicated in Figure A.14. Along route 1 we have $x = y$, so $dl = \sqrt{2} dx$. Thus,

$$\int_P^Q xy^2 dl = \int_0^1 x^3 \sqrt{2} dx = \frac{\sqrt{2}}{4}. \quad (\text{A.70})$$

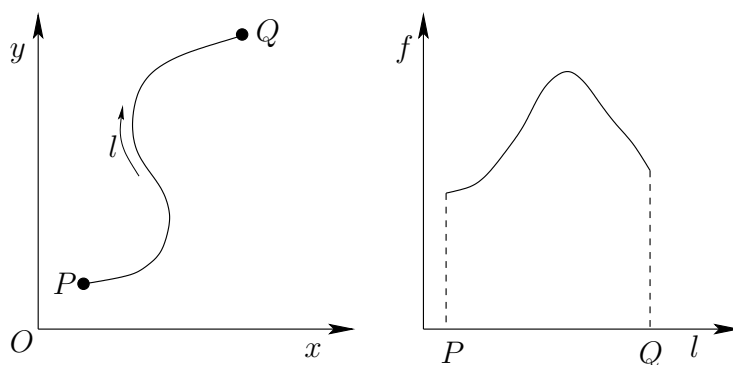


Figure A.13: A line integral.

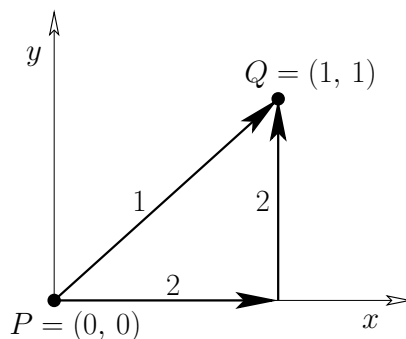


Figure A.14: An example line integral.

The integration along route 2 gives

$$\begin{aligned} \int_P^Q xy^2 dl &= \int_0^1 xy^2 dx \Big|_{y=0} + \int_0^1 xy^2 dy \Big|_{x=1} \\ &= 0 + \int_0^1 y^2 dy = \frac{1}{3}. \end{aligned} \quad (\text{A.71})$$

Note that the integral depends on the route taken between the initial and final points.

The most common type of line integral is one in which the contributions from dx and dy are evaluated separately, rather than through the path length dl ; that is,

$$\int_P^Q [f(x, y) dx + g(x, y) dy]. \quad (\text{A.72})$$

As an example of this, consider the integral

$$\int_P^Q [y dx + x^3 dy] \quad (\text{A.73})$$

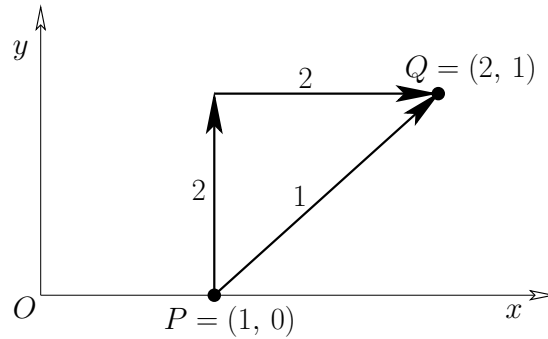


Figure A.15: An example line integral.

along the two routes indicated in Figure A.15. Along route 1 we have $x = y + 1$ and $dx = dy$, so

$$\int_P^Q [y dx + x^3 dy] = \int_0^1 [y dy + (y + 1)^3 dy] = \frac{17}{4}. \tag{A.74}$$

Along route 2,

$$\int_P^Q [y dx + x^3 dy] = \int_0^1 x^3 dy \Big|_{x=1} + \int_1^2 y dx \Big|_{y=1} = \frac{7}{4}. \tag{A.75}$$

Again, the integral depends on the path of integration.

Suppose that we have a line integral that does not depend on the path of integration. It follows that

$$\int_P^Q (f dx + g dy) = F(Q) - F(P) \tag{A.76}$$

for some function F . Given $F(P)$ for one point P in the x - y plane, then

$$F(Q) = F(P) + \int_P^Q (f dx + g dy) \tag{A.77}$$

defines $F(Q)$ for all other points in the plane. We can then draw a contour map of $F(x, y)$. The line integral between points P and Q is simply the change in height in the contour map between these two points:

$$\int_P^Q (f dx + g dy) = \int_P^Q dF(x, y) = F(Q) - F(P). \tag{A.78}$$

Thus,

$$dF(x, y) = f(x, y) dx + g(x, y) dy. \tag{A.79}$$

For instance, if $F = x^3 y$ then $dF = 3x^2 y dx + x^3 dy$ and

$$\int_P^Q (3x^2 y dx + x^3 dy) = [x^3 y]_P^Q \tag{A.80}$$

is independent of the path of integration.

It is clear that there are two distinct types of line integral. Those that depend only on their endpoints and not on the path of integration, and those that depend both on their endpoints and the integration path. Later on, we shall learn how to distinguish between these two types. (See Section A.18.)

A.14 Vector Line Integrals

A *vector field* is defined as a set of vectors associated with each point in space. For instance, the velocity $\mathbf{v}(\mathbf{r})$ in a moving liquid (e.g., a whirlpool) constitutes a vector field. By analogy, a *scalar field* is a set of scalars associated with each point in space. An example of a scalar field is the temperature distribution $T(\mathbf{r})$ in a furnace.

Consider a general vector field $\mathbf{A}(\mathbf{r})$. Let $d\mathbf{r} \equiv (dx, dy, dz)$ be the vector element of line length. Vector line integrals often arise as

$$\int_P^Q \mathbf{A} \cdot d\mathbf{r} = \int_P^Q (A_x dx + A_y dy + A_z dz). \quad (\text{A.81})$$

For instance, if \mathbf{A} is a force-field then the line integral is the work done in going from P to Q .

As an example, consider the work done by a repulsive inverse-square central field, $\mathbf{F} = -\mathbf{r}/|r^3|$. The element of work done is $dW = \mathbf{F} \cdot d\mathbf{r}$. Take $P = (\infty, 0, 0)$ and $Q = (a, 0, 0)$. Route 1 is along the x -axis, so

$$W = \int_{\infty}^a \left(-\frac{1}{x^2} \right) dx = \left[\frac{1}{x} \right]_{\infty}^a = \frac{1}{a}. \quad (\text{A.82})$$

The second route is, firstly, around a large circle ($r = \text{constant}$) to the point $(a, \infty, 0)$, and then parallel to the y -axis. See Figure A.16. In the first part, no work is done, because \mathbf{F} is perpendicular to $d\mathbf{r}$. In the second part,

$$W = \int_{\infty}^0 \frac{-y dy}{(a^2 + y^2)^{3/2}} = \left[\frac{1}{(y^2 + a^2)^{1/2}} \right]_{\infty}^0 = \frac{1}{a}. \quad (\text{A.83})$$

In this case, the integral is independent of the path. However, not all vector line integrals are path independent.

A.15 Surface Integrals

Let us take a surface S , that is not necessarily co-planar, and divide it up into (scalar) elements δS_i . Then

$$\iint_S f(x, y, z) dS = \lim_{\delta S_i \rightarrow 0} \sum_i f(x, y, z) \delta S_i \quad (\text{A.84})$$

is a surface integral. For instance, the volume of water in a lake of depth $D(x, y)$ is

$$V = \iint D(x, y) dS. \quad (\text{A.85})$$

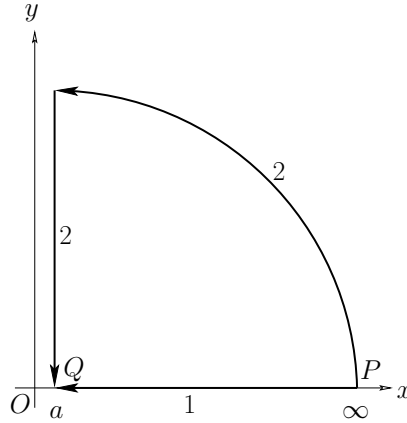


Figure A.16: An example vector line integral.

To evaluate this integral, we must split the calculation into two ordinary integrals. The volume in the strip shown in Figure A.17 is

$$\left[\int_{x_1}^{x_2} D(x, y) dx \right] dy. \tag{A.86}$$

Note that the limits x_1 and x_2 depend on y . The total volume is the sum over all strips: that is,

$$V = \int_{y_1}^{y_2} dy \left[\int_{x_1(y)}^{x_2(y)} D(x, y) dx \right] \equiv \iint_S D(x, y) dx dy. \tag{A.87}$$

Of course, the integral can be evaluated by taking the strips the other way around: that is,

$$V = \int_{x_1}^{x_2} dx \int_{y_1(x)}^{y_2(x)} D(x, y) dy. \tag{A.88}$$

Interchanging the order of integration is a very powerful and useful trick. But great care must be taken when evaluating the limits.

For example, consider

$$\iint_S xy^2 dx dy, \tag{A.89}$$

where S is shown in Figure A.18. Suppose that we evaluate the x integral first:

$$dy \left(\int_0^{1-y} xy^2 dx \right) = y^2 dy \left[\frac{x^2}{2} \right]_0^{1-y} = \frac{y^2}{2} (1-y)^2 dy. \tag{A.90}$$

Let us now evaluate the y integral:

$$\int_0^1 \left(\frac{y^2}{2} - y^3 + \frac{y^4}{2} \right) dy = \frac{1}{60}. \tag{A.91}$$

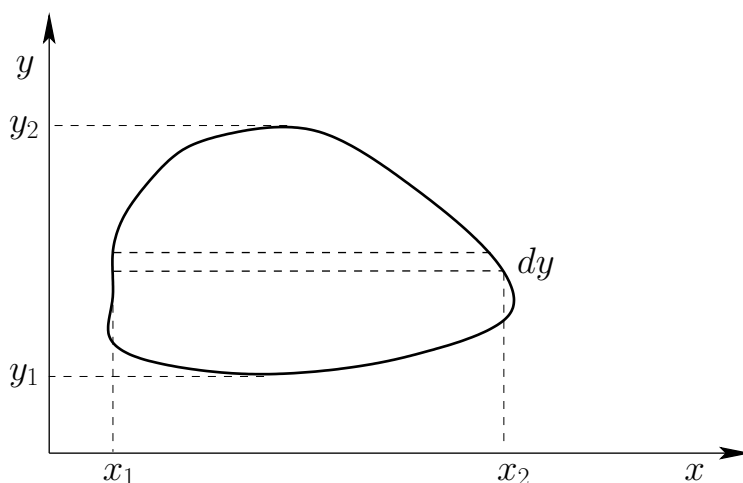


Figure A.17: Decomposition of a surface integral.

We can also evaluate the integral by interchanging the order of integration:

$$\int_0^1 x dx \int_0^{1-x} y^2 dy = \int_0^1 \frac{x}{3} (1-x)^3 dx = \frac{1}{60}. \quad (\text{A.92})$$

In some cases, a surface integral is just the product of two separate integrals. For instance,

$$\int \int_S x^2 y dx dy \quad (\text{A.93})$$

where S is a unit square. This integral can be written

$$\int_0^1 dx \int_0^1 x^2 y dy = \left(\int_0^1 x^2 dx \right) \left(\int_0^1 y dy \right) = \frac{1}{3} \frac{1}{2} = \frac{1}{6}, \quad (\text{A.94})$$

because the limits are both independent of the other variable.

A.16 Vector Surface Integrals

Surface integrals often occur during vector analysis. For instance, the rate of flow of a liquid of velocity \mathbf{v} through an infinitesimal surface of vector area $d\mathbf{S}$ is $\mathbf{v} \cdot d\mathbf{S}$. The net rate of flow through a surface \mathbf{S} made up of very many infinitesimal surfaces is

$$\int \int_S \mathbf{v} \cdot d\mathbf{S} = \lim_{dS \rightarrow 0} \left[\sum v \cos \theta dS \right], \quad (\text{A.95})$$

where θ is the angle subtended between the normal to the surface and the flow velocity.

Analogously to line integrals, most surface integrals depend both on the surface and the rim. But some (very important) integrals depend only on the rim, and not on the nature of the surface

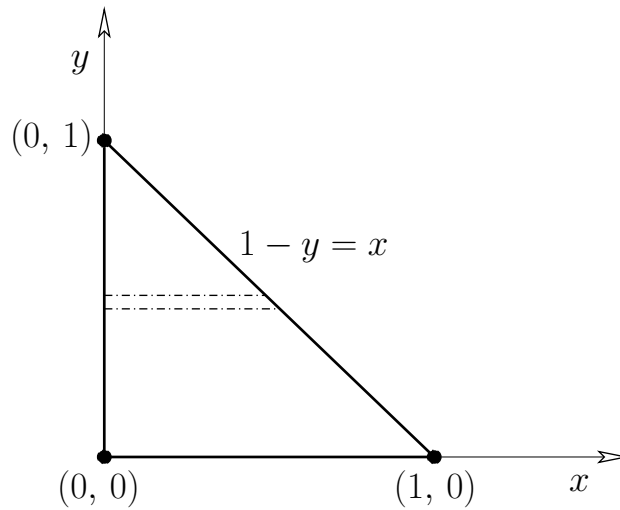


Figure A.18: An example surface integral.

which spans it. As an example of this, consider incompressible fluid flow between two surfaces S_1 and S_2 that end on the same rim. (See Figure A.23.) The volume between the surfaces is constant, so what goes in must come out, and

$$\iint_{S_1} \mathbf{v} \cdot d\mathbf{S} = \iint_{S_2} \mathbf{v} \cdot d\mathbf{S}. \tag{A.96}$$

It follows that

$$\iint \mathbf{v} \cdot d\mathbf{S} \tag{A.97}$$

depends only on the rim, and not on the form of surfaces S_1 and S_2 .

A.17 Volume Integrals

A volume integral takes the form

$$\iiint_V f(x, y, z) dV, \tag{A.98}$$

where V is some volume, and $dV = dx dy dz$ is a small volume element. The volume element is sometimes written $d^3\mathbf{r}$, or even $d\tau$.

As an example of a volume integral, let us evaluate the center of gravity of a solid pyramid. Suppose that the pyramid has a square base of side a , a height a , and is composed of material of uniform density. Let the centroid of the base lie at the origin, and let the apex lie at $(0, 0, a)$. By symmetry, the center of mass lies on the line joining the centroid to the apex. In fact, the height of the center of mass is given by

$$\bar{z} = \iiint z dV / \iiint dV. \tag{A.99}$$

The bottom integral is just the volume of the pyramid, and can be written

$$\begin{aligned}\iiint dV &= \int_0^a dz \int_{-(a-z)/2}^{(a-z)/2} dy \int_{-(a-z)/2}^{(a-z)/2} dx = \int_0^a (a-z)^2 dz = \int_0^a (a^2 - 2az + z^2) dz \\ &= [a^2 z - az^2 + z^3/3]_0^a = \frac{1}{3} a^3.\end{aligned}\quad (\text{A.100})$$

Here, we have evaluated the z -integral last because the limits of the x - and y - integrals are z -dependent. The top integral takes the form

$$\begin{aligned}\iiint z dV &= \int_0^a z dz \int_{-(a-z)/2}^{(a-z)/2} dy \int_{-(a-z)/2}^{(a-z)/2} dx = \int_0^a z(a-z)^2 dz = \int_0^a (za^2 - 2az^2 + z^3) dz \\ &= [a^2 z^2/2 - 2az^3/3 + z^4/4]_0^a = \frac{1}{12} a^4.\end{aligned}\quad (\text{A.101})$$

Thus,

$$\bar{z} = \frac{1}{12} a^4 \bigg/ \frac{1}{3} a^3 = \frac{1}{4} a. \quad (\text{A.102})$$

In other words, the center of mass of a pyramid lies one quarter of the way between the centroid of the base and the apex.

A.18 Gradient

A one-dimensional function $f(x)$ has a gradient df/dx that is defined as the slope of the tangent to the curve at x . We wish to extend this idea to cover scalar fields in two and three dimensions.

Consider a two-dimensional scalar field $h(x, y)$, which is (say) height above sea-level in a hilly region. Let $d\mathbf{r} \equiv (dx, dy)$ be an element of horizontal distance. Consider dh/dr , where dh is the change in height after moving an infinitesimal distance $d\mathbf{r}$. This quantity is somewhat like the one-dimensional gradient, except that dh depends on the direction of $d\mathbf{r}$, as well as its magnitude. In the immediate vicinity of some point P , the slope reduces to an inclined plane. See Figure A.19. The largest value of dh/dr is straight up the slope. It is easily shown that for any other direction

$$\frac{dh}{dr} = \left(\frac{dh}{dr} \right)_{\max} \cos \theta, \quad (\text{A.103})$$

where θ is the angle shown in Figure A.19. Let us define a two-dimensional vector, **grad** h , called the *gradient* of h , whose magnitude is $(dh/dr)_{\max}$, and whose direction is the direction of steepest ascent. The $\cos \theta$ variation exhibited in the previous expression ensures that the component of **grad** h in any direction is equal to dh/dr for that direction.

The component of dh/dr in the x -direction can be obtained by plotting out the profile of h at constant y , and then finding the slope of the tangent to the curve at given x . This quantity is known as the *partial derivative* of h with respect to x at constant y , and is denoted $(\partial h/\partial x)_y$. Likewise, the

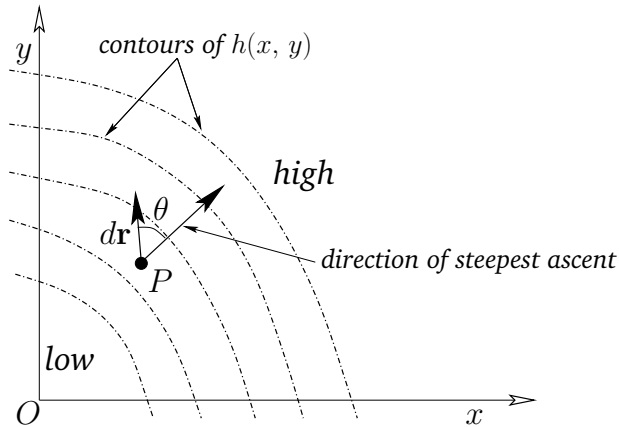


Figure A.19: A two-dimensional gradient.

gradient of the profile at constant x is written $(\partial h/\partial y)_x$. Note that the subscripts denoting constant- x and constant- y are usually omitted, unless there is any ambiguity. It follows that in component form

$$\mathbf{grad} h \equiv \left(\frac{\partial h}{\partial x}, \frac{\partial h}{\partial y} \right). \quad (\text{A.104})$$

The equation of the tangent plane at $P = (x_0, y_0)$ is

$$h_T(x, y) = h(x_0, y_0) + \alpha(x - x_0) + \beta(y - y_0). \quad (\text{A.105})$$

This plane has the same local gradients as $h(x, y)$, so

$$\alpha = \frac{\partial h}{\partial x}, \quad \beta = \frac{\partial h}{\partial y}, \quad (\text{A.106})$$

by differentiation of the previous equation. For small $dx = x - x_0$ and $dy = y - y_0$, the function h is coincident with the tangent plane, so

$$dh = \frac{\partial h}{\partial x} dx + \frac{\partial h}{\partial y} dy. \quad (\text{A.107})$$

But, $\mathbf{grad} h \equiv (\partial h/\partial x, \partial h/\partial y)$ and $d\mathbf{r} \equiv (dx, dy)$, so

$$dh = \mathbf{grad} h \cdot d\mathbf{r}. \quad (\text{A.108})$$

Incidentally, the previous equation demonstrates that $\mathbf{grad} h$ is a proper vector, because the left-hand side is a scalar, and, according to the properties of the dot product, the right-hand side is also a scalar provided that $d\mathbf{r}$ and $\mathbf{grad} h$ are both proper vectors ($d\mathbf{r}$ is an obvious vector, because it is directly derived from displacements).

Consider, now, a three-dimensional temperature distribution $T(x, y, z)$ in (say) a reaction vessel. Let us define $\mathbf{grad} T$, as before, as a vector whose magnitude is $(dT/dr)_{\max}$, and whose direction is the direction of the maximum gradient. This vector is written in component form

$$\mathbf{grad} T \equiv \left(\frac{\partial T}{\partial x}, \frac{\partial T}{\partial y}, \frac{\partial T}{\partial z} \right). \quad (\text{A.109})$$

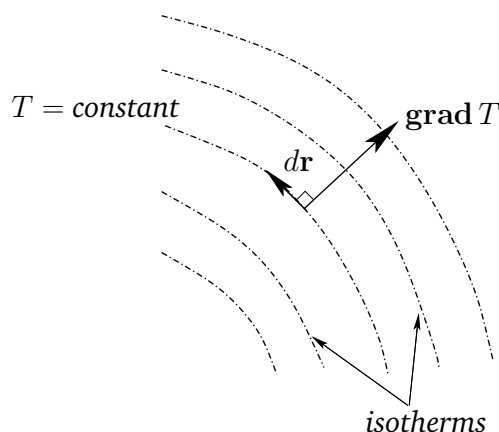


Figure A.20: Isotherms.

Here, $\partial T/\partial x \equiv (\partial T/\partial x)_{y,z}$ is the gradient of the one-dimensional temperature profile at constant y and z . The change in T in going from point P to a neighbouring point offset by $d\mathbf{r} \equiv (dx, dy, dz)$ is

$$dT = \frac{\partial T}{\partial x} dx + \frac{\partial T}{\partial y} dy + \frac{\partial T}{\partial z} dz. \quad (\text{A.110})$$

In vector form, this becomes

$$dT = \mathbf{grad} T \cdot d\mathbf{r}. \quad (\text{A.111})$$

Suppose that $dT = 0$ for some $d\mathbf{r}$. It follows that

$$dT = \mathbf{grad} T \cdot d\mathbf{r} = 0. \quad (\text{A.112})$$

So, $d\mathbf{r}$ is perpendicular to $\mathbf{grad} T$. Because $dT = 0$ along so-called “isotherms” (i.e., contours of the temperature), we conclude that the isotherms (contours) are everywhere perpendicular to $\mathbf{grad} T$. See Figure A.20.

It is, of course, possible to integrate dT . For instance, the line integral of dT between points P and Q is written

$$\int_P^Q dT = \int_P^Q \mathbf{grad} T \cdot d\mathbf{r} = T(Q) - T(P). \quad (\text{A.113})$$

This integral is clearly independent of the path taken between P and Q , so $\int_P^Q \mathbf{grad} T \cdot d\mathbf{r}$ must be path independent.

Consider a vector field $\mathbf{A}(\mathbf{r})$. In general, the line integral $\int_P^Q \mathbf{A} \cdot d\mathbf{r}$ depends on the path taken between the end points, but for some special vector fields the integral is path independent. Such fields are called *conservative* fields. It can be shown that if \mathbf{A} is a conservative field then $\mathbf{A} = \mathbf{grad} V$ for some scalar field V . The proof of this is straightforward. Keeping P fixed, we have

$$\int_P^Q \mathbf{A} \cdot d\mathbf{r} = V(Q), \quad (\text{A.114})$$

where $V(Q)$ is a well-defined function, due to the path-independent nature of the line integral. Consider moving the position of the end point by an infinitesimal amount dx in the x -direction. We have

$$V(Q + dx) = V(Q) + \int_Q^{Q+dx} \mathbf{A} \cdot d\mathbf{r} = V(Q) + A_x dx. \quad (\text{A.115})$$

Hence,

$$\frac{\partial V}{\partial x} = A_x, \quad (\text{A.116})$$

with analogous relations for the other components of \mathbf{A} . It follows that

$$\mathbf{A} = \mathbf{grad} V. \quad (\text{A.117})$$

In classical dynamics, the force due to gravity is a good example of a conservative field. Now, if $\mathbf{A}(\mathbf{r})$ is a force-field then $\int \mathbf{A} \cdot d\mathbf{r}$ is the work done in traversing some path. If \mathbf{A} is conservative then

$$\oint \mathbf{A} \cdot d\mathbf{r} = 0, \quad (\text{A.118})$$

where \oint corresponds to the line integral around a closed loop. The fact that zero net work is done in going around a closed loop is equivalent to the conservation of energy (which is why conservative fields are called “conservative”). A good example of a non-conservative field is the force due to friction. Clearly, a frictional system loses energy in going around a closed cycle, so $\oint \mathbf{A} \cdot d\mathbf{r} \neq 0$.

A.19 Grad Operator

It is useful to define the vector operator

$$\nabla \equiv \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z} \right), \quad (\text{A.119})$$

which is usually called the *grad* or *del* operator. This operator acts on everything to its right in an expression, until the end of the expression or a closing bracket is reached. For instance,

$$\mathbf{grad} f = \nabla f \equiv \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} \right). \quad (\text{A.120})$$

For two scalar fields ϕ and ψ ,

$$\mathbf{grad} (\phi \psi) = \phi \mathbf{grad} \psi + \psi \mathbf{grad} \phi \quad (\text{A.121})$$

can be written more succinctly as

$$\nabla(\phi \psi) = \phi \nabla \psi + \psi \nabla \phi. \quad (\text{A.122})$$

Suppose that we rotate the coordinate axes through an angle θ about Oz . By analogy with Equations (A.17)–(A.19), the old coordinates (x, y, z) are related to the new ones (x', y', z') via

$$x = x' \cos \theta - y' \sin \theta, \quad (\text{A.123})$$

$$y = x' \sin \theta + y' \cos \theta, \quad (\text{A.124})$$

$$z = z'. \quad (\text{A.125})$$

Now,

$$\frac{\partial}{\partial x'} = \left(\frac{\partial x}{\partial x'} \right)_{y', z'} \frac{\partial}{\partial x} + \left(\frac{\partial y}{\partial x'} \right)_{y', z'} \frac{\partial}{\partial y} + \left(\frac{\partial z}{\partial x'} \right)_{y', z'} \frac{\partial}{\partial z}, \quad (\text{A.126})$$

giving

$$\frac{\partial}{\partial x'} = \cos \theta \frac{\partial}{\partial x} + \sin \theta \frac{\partial}{\partial y}, \quad (\text{A.127})$$

and

$$\nabla_{x'} = \cos \theta \nabla_x + \sin \theta \nabla_y. \quad (\text{A.128})$$

It can be seen, from Equations (A.20)–(A.22), that the differential operator ∇ transforms in an analogous manner to a vector. This is another proof that ∇f is a good vector.

A.20 Divergence

Let us start with a vector field $\mathbf{A}(\mathbf{r})$. Consider $\oint_S \mathbf{A} \cdot d\mathbf{S}$ over some closed surface S , where $d\mathbf{S}$ denotes an outward pointing surface element. This surface integral is usually called the *flux* of \mathbf{A} out of S . If \mathbf{A} represents the velocity of some fluid then $\oint_S \mathbf{A} \cdot d\mathbf{S}$ is the rate of fluid flow out of S .

If \mathbf{A} is constant in space then it is easily demonstrated that the net flux out of S is zero. In fact,

$$\oint \mathbf{A} \cdot d\mathbf{S} = \mathbf{A} \cdot \oint d\mathbf{S} = \mathbf{A} \cdot \mathbf{S} = 0, \quad (\text{A.129})$$

because the vector area \mathbf{S} of a closed surface is zero.

Suppose, now, that \mathbf{A} is not uniform in space. Consider a very small rectangular volume over which \mathbf{A} hardly varies. The contribution to $\oint \mathbf{A} \cdot d\mathbf{S}$ from the two faces normal to the x -axis is

$$A_x(x + dx) dy dz - A_x(x) dy dz = \frac{\partial A_x}{\partial x} dx dy dz = \frac{\partial A_x}{\partial x} dV, \quad (\text{A.130})$$

where $dV = dx dy dz$ is the volume element. (See Figure A.21.) There are analogous contributions from the sides normal to the y - and z -axes, so the total of all the contributions is

$$\oint \mathbf{A} \cdot d\mathbf{S} = \left(\frac{\partial A_x}{\partial x} + \frac{\partial A_y}{\partial y} + \frac{\partial A_z}{\partial z} \right) dV. \quad (\text{A.131})$$

The *divergence* of a vector field is defined

$$\text{div } \mathbf{A} = \nabla \cdot \mathbf{A} = \frac{\partial A_x}{\partial x} + \frac{\partial A_y}{\partial y} + \frac{\partial A_z}{\partial z}. \quad (\text{A.132})$$

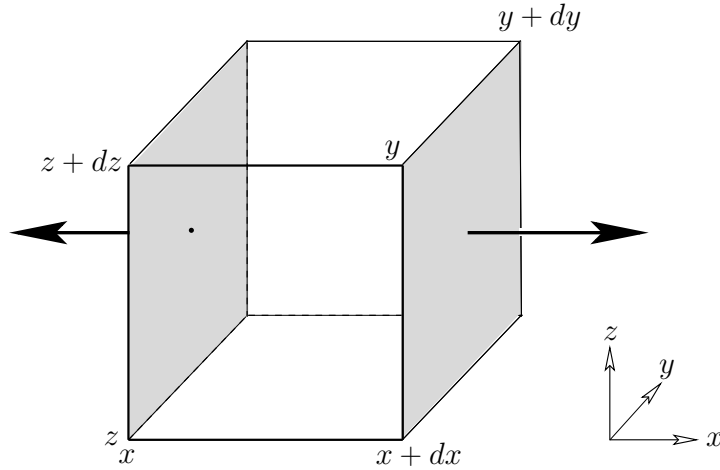


Figure A.21: Flux of a vector field out of a small box.

Divergence is a good scalar (i.e., it is coordinate independent), because it is the dot product of the vector operator ∇ with \mathbf{A} . The formal definition of $\nabla \cdot \mathbf{A}$ is

$$\nabla \cdot \mathbf{A} = \lim_{dV \rightarrow 0} \frac{\oint \mathbf{A} \cdot d\mathbf{S}}{dV}. \quad (\text{A.133})$$

This definition is independent of the shape of the infinitesimal volume element.

One of the most important results in vector field theory is the so-called *divergence theorem*. This states that for any volume V surrounded by a closed surface S ,

$$\oint_S \mathbf{A} \cdot d\mathbf{S} = \int_V \nabla \cdot \mathbf{A} \, dV, \quad (\text{A.134})$$

where $d\mathbf{S}$ is an outward pointing volume element. The proof is very straightforward. We divide up the volume into very many infinitesimal cubes, and sum $\int \mathbf{A} \cdot d\mathbf{S}$ over all of the surfaces. The contributions from the interior surfaces cancel out, leaving just the contribution from the outer surface. (See Figure A.22.) We can use Equation (A.131) for each cube individually. This tells us that the summation is equivalent to $\int \nabla \cdot \mathbf{A} \, dV$ over the whole volume. Thus, the integral of $\mathbf{A} \cdot d\mathbf{S}$ over the outer surface is equal to the integral of $\nabla \cdot \mathbf{A}$ over the whole volume, which proves the divergence theorem.

Now, for a vector field with $\nabla \cdot \mathbf{A} = 0$,

$$\oint_S \mathbf{A} \cdot d\mathbf{S} = 0 \quad (\text{A.135})$$

for any closed surface S . So, for two surfaces, S_1 and S_2 , on the same rim,

$$\int_{S_1} \mathbf{A} \cdot d\mathbf{S} = \int_{S_2} \mathbf{A} \cdot d\mathbf{S}, \quad (\text{A.136})$$

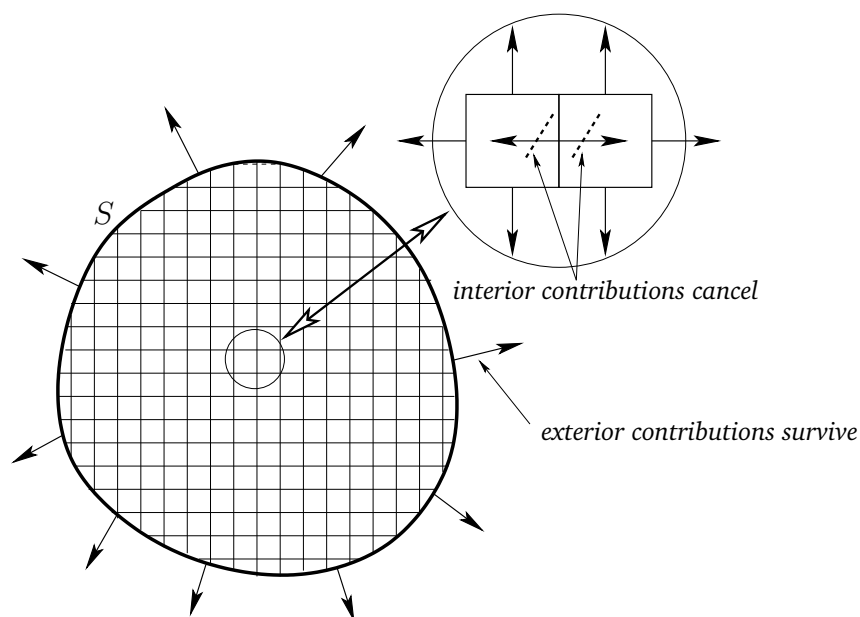


Figure A.22: The divergence theorem.

as illustrated in Figure A.23. (Note that the direction of the surface elements on S_1 has been reversed relative to those on the closed surface. Hence, the sign of the associated surface integral is also reversed.) Thus, if $\nabla \cdot \mathbf{A} = 0$ then the surface integral depends on the rim, but not on the nature of the surface that spans it. On the other hand, if $\nabla \cdot \mathbf{A} \neq 0$ then the integral depends on both the rim and the surface.

Consider an incompressible fluid whose velocity field is \mathbf{v} . It is clear that $\oint \mathbf{v} \cdot d\mathbf{S} = 0$ for any closed surface, because what flows into the surface must flow out again. Thus, according to the divergence theorem, $\int \nabla \cdot \mathbf{v} dV = 0$ for any volume. The only way in which this is possible is if $\nabla \cdot \mathbf{v}$ is everywhere zero. Thus, the velocity components of an incompressible fluid satisfy the

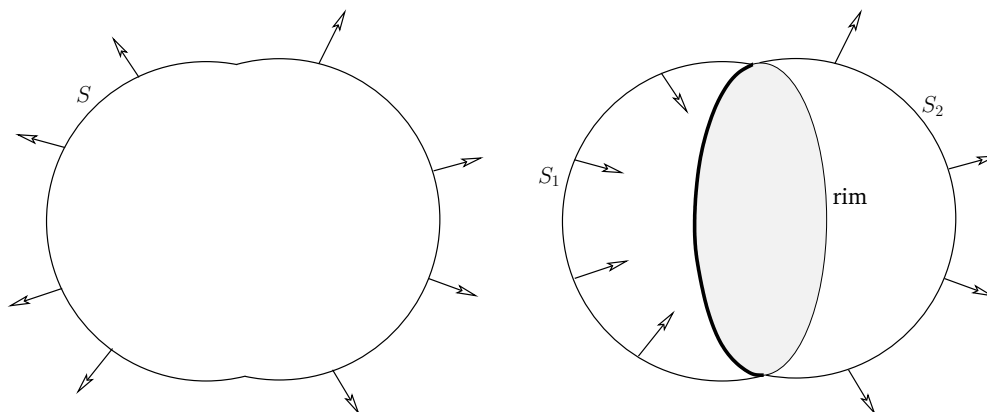


Figure A.23: Two surfaces spanning the same rim (right), and the equivalent closed surface (left).

following differential relation:

$$\frac{\partial v_x}{\partial x} + \frac{\partial v_y}{\partial y} + \frac{\partial v_z}{\partial z} = 0. \quad (\text{A.137})$$

It is sometimes helpful to represent a vector field \mathbf{A} by *lines of force* or *field-lines*. The direction of a line of force at any point is the same as the local direction of \mathbf{A} . The density of lines (i.e., the number of lines crossing a unit surface perpendicular to \mathbf{A}) is equal to $|\mathbf{A}|$. For instance, in Figure A.24, $|\mathbf{A}|$ is larger at point 1 than at point 2. The number of lines crossing a surface element $d\mathbf{S}$ is $\mathbf{A} \cdot d\mathbf{S}$. So, the net number of lines leaving a closed surface is

$$\oint_S \mathbf{A} \cdot d\mathbf{S} = \int_V \nabla \cdot \mathbf{A} \, dV. \quad (\text{A.138})$$

If $\nabla \cdot \mathbf{A} = 0$ then there is no net flux of lines out of any surface. Such a field is called a *solenoidal* vector field. The simplest example of a solenoidal vector field is one in which the lines of force all form closed loops.

A.21 Laplacian Operator

So far we have encountered

$$\nabla\phi = \left(\frac{\partial\phi}{\partial x}, \frac{\partial\phi}{\partial y}, \frac{\partial\phi}{\partial z} \right), \quad (\text{A.139})$$

which is a vector field formed from a scalar field, and

$$\nabla \cdot \mathbf{A} = \frac{\partial A_x}{\partial x} + \frac{\partial A_y}{\partial y} + \frac{\partial A_z}{\partial z}, \quad (\text{A.140})$$

which is a scalar field formed from a vector field. There are two ways in which we can combine gradient and divergence. We can either form the vector field $\nabla(\nabla \cdot \mathbf{A})$ or the scalar field $\nabla \cdot (\nabla\phi)$. The former is not particularly interesting, but the scalar field $\nabla \cdot (\nabla\phi)$ turns up in a great many physical problems, and is, therefore, worthy of discussion.

Let us introduce the heat flow vector \mathbf{h} , which is the rate of flow of heat energy per unit area across a surface perpendicular to the direction of \mathbf{h} . In many substances, heat flows directly down

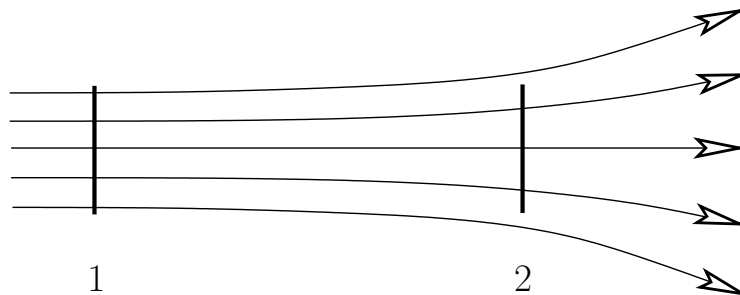


Figure A.24: Divergent lines of force.

the temperature gradient, so that we can write

$$\mathbf{h} = -\kappa \nabla T, \quad (\text{A.141})$$

where κ is the *thermal conductivity*. The net rate of heat flow $\oint_S \mathbf{h} \cdot d\mathbf{S}$ out of some closed surface S must be equal to the rate of decrease of heat energy in the volume V enclosed by S . Thus, we have

$$\oint_S \mathbf{h} \cdot d\mathbf{S} = -\frac{\partial}{\partial t} \left(\int c T dV \right), \quad (\text{A.142})$$

where c is the *specific heat*. It follows from the divergence theorem that

$$\nabla \cdot \mathbf{h} = -c \frac{\partial T}{\partial t}. \quad (\text{A.143})$$

Taking the divergence of both sides of Equation (A.141), and making use of Equation (A.143), we obtain

$$\nabla \cdot (\kappa \nabla T) = c \frac{\partial T}{\partial t}. \quad (\text{A.144})$$

If κ is constant then the previous equation can be written

$$\nabla \cdot (\nabla T) = \frac{c}{\kappa} \frac{\partial T}{\partial t}. \quad (\text{A.145})$$

The scalar field $\nabla \cdot (\nabla T)$ takes the form

$$\begin{aligned} \nabla \cdot (\nabla T) &= \frac{\partial}{\partial x} \left(\frac{\partial T}{\partial x} \right) + \frac{\partial}{\partial y} \left(\frac{\partial T}{\partial y} \right) + \frac{\partial}{\partial z} \left(\frac{\partial T}{\partial z} \right) \\ &= \frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} + \frac{\partial^2 T}{\partial z^2} \equiv \nabla^2 T. \end{aligned} \quad (\text{A.146})$$

Here, the scalar differential operator

$$\nabla^2 \equiv \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} \quad (\text{A.147})$$

is called the *Laplacian*. The Laplacian is a good scalar operator (i.e., it is coordinate independent) because it is formed from a combination of divergence (another good scalar operator) and gradient (a good vector operator).

What is the physical significance of the Laplacian? In one dimension, $\nabla^2 T$ reduces to $\partial^2 T / \partial x^2$. Now, $\partial^2 T / \partial x^2$ is positive if $T(x)$ is concave (from above), and negative if it is convex. So, if T is less than the average of T in its surroundings then $\nabla^2 T$ is positive, and vice versa.

In two dimensions,

$$\nabla^2 T = \frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2}. \quad (\text{A.148})$$

Consider a local minimum of the temperature. At the minimum, the slope of T increases in all directions, so $\nabla^2 T$ is positive. Likewise, $\nabla^2 T$ is negative at a local maximum. Consider, now, a

steep-sided valley in T . Suppose that the bottom of the valley runs parallel to the x -axis. At the bottom of the valley $\partial^2 T / \partial y^2$ is large and positive, whereas $\partial^2 T / \partial x^2$ is small and may even be negative. Thus, $\nabla^2 T$ is positive, and this is associated with T being less than the average local value.

Let us now return to the heat conduction problem:

$$\nabla^2 T = \frac{c}{\kappa} \frac{\partial T}{\partial t}. \quad (\text{A.149})$$

It is clear that if $\nabla^2 T$ is positive in some small region then the value of T there is less than the local average value, so $\partial T / \partial t > 0$: that is, the region heats up. Likewise, if $\nabla^2 T$ is negative then the value of T is greater than the local average value, and heat flows out of the region: that is, $\partial T / \partial t < 0$. Thus, the previous heat conduction equation makes physical sense.

A.22 Curl

Consider a vector field $\mathbf{A}(\mathbf{r})$, and a loop that lies in one plane. The integral of \mathbf{A} around this loop is written $\oint \mathbf{A} \cdot d\mathbf{r}$, where $d\mathbf{r}$ is a line element of the loop. If \mathbf{A} is a conservative field then $\mathbf{A} = \nabla\phi$ and $\oint \mathbf{A} \cdot d\mathbf{r} = 0$ for all loops. In general, for a non-conservative field, $\oint \mathbf{A} \cdot d\mathbf{r} \neq 0$.

For a small loop, we expect $\oint \mathbf{A} \cdot d\mathbf{r}$ to be proportional to the area of the loop. Moreover, for a fixed-area loop, we expect $\oint \mathbf{A} \cdot d\mathbf{r}$ to depend on the orientation of the loop. One particular orientation will give the maximum value: $\oint \mathbf{A} \cdot d\mathbf{r} = I_{\max}$. If the loop subtends an angle θ with this optimum orientation then we expect $I = I_{\max} \cos \theta$. Let us introduce the vector field **curl** \mathbf{A} whose magnitude is

$$|\mathbf{curl} \mathbf{A}| = \lim_{dS \rightarrow 0} \frac{\oint \mathbf{A} \cdot d\mathbf{r}}{dS} \quad (\text{A.150})$$

for the orientation giving I_{\max} . Here, dS is the area of the loop. The direction of **curl** \mathbf{A} is perpendicular to the plane of the loop, when it is in the orientation giving I_{\max} , with the sense given by a right-hand circulation rule.

Let us now express **curl** \mathbf{A} in terms of the components of \mathbf{A} . First, we shall evaluate $\oint \mathbf{A} \cdot d\mathbf{r}$ around a small rectangle in the y - z plane, as shown in Figure A.25. The contribution from sides 1 and 3 is

$$A_z(y + dy) dz - A_z(y) dz = \frac{\partial A_z}{\partial y} dy dz. \quad (\text{A.151})$$

The contribution from sides 2 and 4 is

$$-A_y(z + dz) dy + A_y(z) dy = -\frac{\partial A_y}{\partial z} dy dz. \quad (\text{A.152})$$

So, the total of all contributions gives

$$\oint \mathbf{A} \cdot d\mathbf{r} = \left(\frac{\partial A_z}{\partial y} - \frac{\partial A_y}{\partial z} \right) dS, \quad (\text{A.153})$$

where $dS = dy dz$ is the area of the loop.

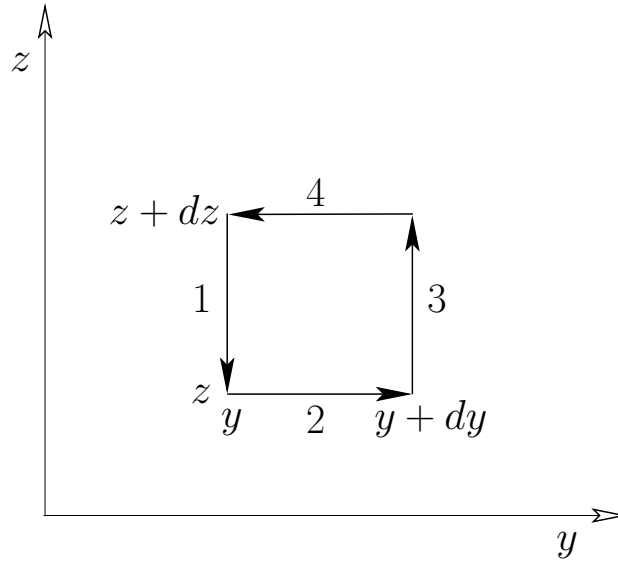


Figure A.25: A vector line integral around a small rectangular loop in the y - z plane.

Consider a non-rectangular (but still small) loop in the y - z plane. We can divide it into rectangular elements, and form $\oint \mathbf{A} \cdot d\mathbf{r}$ over all the resultant loops. The interior contributions cancel, so we are just left with the contribution from the outer loop. Also, the area of the outer loop is the sum of all the areas of the inner loops. We conclude that

$$\oint \mathbf{A} \cdot d\mathbf{r} = \left(\frac{\partial A_z}{\partial y} - \frac{\partial A_y}{\partial z} \right) dS_x \quad (\text{A.154})$$

is valid for a small loop $d\mathbf{S} = (dS_x, 0, 0)$ of any shape in the y - z plane. Likewise, we can show that if the loop is in the x - z plane then $d\mathbf{S} = (0, dS_y, 0)$ and

$$\oint \mathbf{A} \cdot d\mathbf{r} = \left(\frac{\partial A_x}{\partial z} - \frac{\partial A_z}{\partial x} \right) dS_y. \quad (\text{A.155})$$

Finally, if the loop is in the x - y plane then $d\mathbf{S} = (0, 0, dS_z)$ and

$$\oint \mathbf{A} \cdot d\mathbf{r} = \left(\frac{\partial A_y}{\partial x} - \frac{\partial A_x}{\partial y} \right) dS_z. \quad (\text{A.156})$$

Imagine an arbitrary loop of vector area $d\mathbf{S} = (dS_x, dS_y, dS_z)$. We can construct this out of three vector areas, 1, 2, and 3, directed in the x -, y -, and z -directions, respectively, as indicated in Figure A.26. If we form the line integral around all three loops then the interior contributions cancel, and we are left with the line integral around the original loop. Thus,

$$\oint \mathbf{A} \cdot d\mathbf{r} = \oint \mathbf{A} \cdot d\mathbf{r}_1 + \oint \mathbf{A} \cdot d\mathbf{r}_2 + \oint \mathbf{A} \cdot d\mathbf{r}_3, \quad (\text{A.157})$$

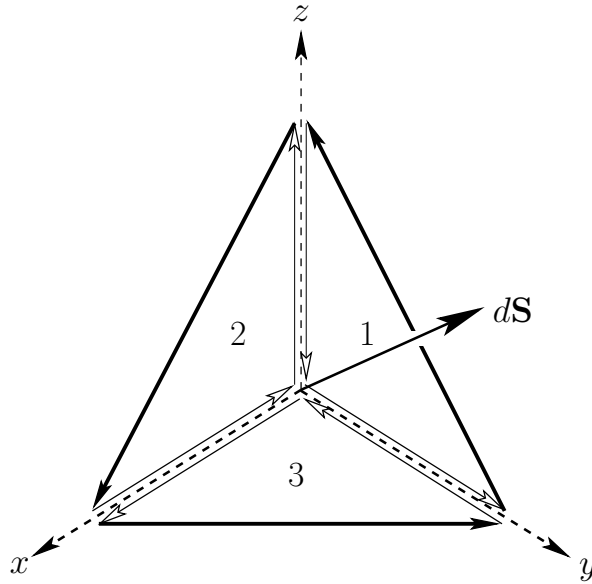


Figure A.26: Decomposition of a vector area into its Cartesian components.

giving

$$\oint \mathbf{A} \cdot d\mathbf{r} = \mathbf{curl} \mathbf{A} \cdot d\mathbf{S} = |\mathbf{curl} \mathbf{A}| |d\mathbf{S}| \cos \theta, \quad (\text{A.158})$$

where

$$\mathbf{curl} \mathbf{A} = \left(\frac{\partial A_z}{\partial y} - \frac{\partial A_y}{\partial z}, \frac{\partial A_x}{\partial z} - \frac{\partial A_z}{\partial x}, \frac{\partial A_y}{\partial x} - \frac{\partial A_x}{\partial y} \right), \quad (\text{A.159})$$

and θ is the angle subtended between the directions of $\mathbf{curl} \mathbf{A}$ and $d\mathbf{S}$. Note that

$$\mathbf{curl} \mathbf{A} = \nabla \times \mathbf{A}. \quad (\text{A.160})$$

This demonstrates that $\nabla \times \mathbf{A}$ is a good vector field, because it is the cross product of the ∇ operator (a good vector operator) and the vector field \mathbf{A} .

Consider a solid body rotating about the z -axis. The angular velocity is given by $\boldsymbol{\omega} = (0, 0, \omega)$, so the rotation velocity at position \mathbf{r} is

$$\mathbf{v} = \boldsymbol{\omega} \times \mathbf{r}. \quad (\text{A.161})$$

[See Equation (A.52).] Let us evaluate $\nabla \times \mathbf{v}$ on the axis of rotation. The x -component is proportional to the integral $\oint \mathbf{v} \cdot d\mathbf{r}$ around a loop in the y - z plane. This is plainly zero. Likewise, the y -component is also zero. The z -component is $\oint \mathbf{v} \cdot d\mathbf{r}/dS$ around some loop in the x - y plane. Consider a circular loop. We have $\oint \mathbf{v} \cdot d\mathbf{r} = 2\pi r \omega r$ with $dS = \pi r^2$. Here, r is the perpendicular distance from the rotation axis. It follows that $(\nabla \times \mathbf{v})_z = 2\omega$, which is independent of r . So, on the axis, $\nabla \times \mathbf{v} = (0, 0, 2\omega)$. Off the axis, at position \mathbf{r}_0 , we can write

$$\mathbf{v} = \boldsymbol{\omega} \times (\mathbf{r} - \mathbf{r}_0) + \boldsymbol{\omega} \times \mathbf{r}_0. \quad (\text{A.162})$$

The first part has the same curl as the velocity field on the axis, and the second part has zero curl, because it is constant. Thus, $\nabla \times \mathbf{v} = (0, 0, 2\omega)$ everywhere in the body. This allows us to form a physical picture of $\nabla \times \mathbf{A}$. If we imagine $\mathbf{A}(\mathbf{r})$ as the velocity field of some fluid then $\nabla \times \mathbf{A}$ at any given point is equal to twice the local angular rotation velocity: that is, 2ω . Hence, a vector field with $\nabla \times \mathbf{A} = \mathbf{0}$ everywhere is said to be *irrotational*.

Another important result of vector field theory is the *curl theorem*:

$$\oint_C \mathbf{A} \cdot d\mathbf{r} = \int_S \nabla \times \mathbf{A} \cdot d\mathbf{S}, \quad (\text{A.163})$$

for some (non-planar) surface S bounded by a rim C . This theorem can easily be proved by splitting the loop up into many small rectangular loops, and forming the integral around all of the resultant loops. All of the contributions from the interior loops cancel, leaving just the contribution from the outer rim. Making use of Equation (A.158) for each of the small loops, we can see that the contribution from all of the loops is also equal to the integral of $\nabla \times \mathbf{A} \cdot d\mathbf{S}$ across the whole surface. This proves the theorem.

One immediate consequence of the curl theorem is that $\nabla \times \mathbf{A}$ is “incompressible.” Consider any two surfaces, S_1 and S_2 , that share the same rim. (See Figure A.23.) It is clear from the curl theorem that $\int \nabla \times \mathbf{A} \cdot d\mathbf{S}$ is the same for both surfaces. Thus, it follows that $\oint \nabla \times \mathbf{A} \cdot d\mathbf{S} = 0$ for any closed surface. However, we have from the divergence theorem that $\oint \nabla \times \mathbf{A} \cdot d\mathbf{S} = \int \nabla \cdot (\nabla \times \mathbf{A}) dV = 0$ for any volume. Hence,

$$\nabla \cdot (\nabla \times \mathbf{A}) \equiv 0. \quad (\text{A.164})$$

So, $\nabla \times \mathbf{A}$ is a solenoidal field.

We have seen that for a conservative field $\oint \mathbf{A} \cdot d\mathbf{r} = 0$ for any loop. This is entirely equivalent to $\mathbf{A} = \nabla\phi$. However, the magnitude of $\nabla \times \mathbf{A}$ is $\lim_{dS \rightarrow 0} \oint \mathbf{A} \cdot d\mathbf{r}/dS$ for some particular loop. It is clear then that $\nabla \times \mathbf{A} = \mathbf{0}$ for a conservative field. In other words,

$$\nabla \times (\nabla\phi) \equiv \mathbf{0}. \quad (\text{A.165})$$

Thus, a conservative field is also an irrotational one.

A.23 Curvilinear Coordinates

In the cylindrical coordinate system, the Cartesian coordinates x and y are replaced by $r = (x^2 + y^2)^{1/2}$ and $\theta = \tan^{-1}(y/x)$. Here, r is the perpendicular distance from the z -axis, and θ the angle subtended between the perpendicular radius vector and the x -axis. See Figure A.27. A general vector \mathbf{A} is thus written

$$\mathbf{A} = A_r \mathbf{e}_r + A_\theta \mathbf{e}_\theta + A_z \mathbf{e}_z, \quad (\text{A.166})$$

where $\mathbf{e}_r = \nabla r/|\nabla r|$ and $\mathbf{e}_\theta = \nabla\theta/|\nabla\theta|$. See Figure A.27. Note that the unit vectors \mathbf{e}_r , \mathbf{e}_θ , and \mathbf{e}_z are mutually orthogonal. Hence, $A_r = \mathbf{A} \cdot \mathbf{e}_r$, et cetera. The volume element in this coordinate system is $d^3\mathbf{r} = r dr d\theta dz$. Moreover, gradient, divergence, and curl take the forms

$$\nabla V = \frac{\partial V}{\partial r} \mathbf{e}_r + \frac{1}{r} \frac{\partial V}{\partial \theta} \mathbf{e}_\theta + \frac{\partial V}{\partial z} \mathbf{e}_z, \quad (\text{A.167})$$

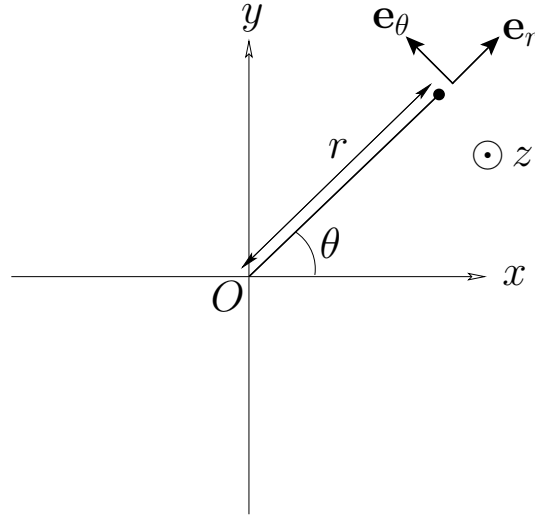


Figure A.27: Cylindrical coordinates.

$$\nabla \cdot \mathbf{A} = \frac{1}{r} \frac{\partial}{\partial r} (r A_r) + \frac{1}{r} \frac{\partial A_\theta}{\partial \theta} + \frac{\partial A_z}{\partial z}, \quad (\text{A.168})$$

$$\begin{aligned} \nabla \times \mathbf{A} = & \left(\frac{1}{r} \frac{\partial A_z}{\partial \theta} - \frac{\partial A_\theta}{\partial z} \right) \mathbf{e}_r + \left(\frac{\partial A_r}{\partial z} - \frac{\partial A_z}{\partial r} \right) \mathbf{e}_\theta \\ & + \left(\frac{1}{r} \frac{\partial}{\partial r} (r A_\theta) - \frac{1}{r} \frac{\partial A_r}{\partial \theta} \right) \mathbf{e}_z, \end{aligned} \quad (\text{A.169})$$

respectively. Here, $V(\mathbf{r})$ is a general vector field, and $\mathbf{A}(\mathbf{r})$ a general scalar field. Finally, the Laplacian is written

$$\nabla^2 V = \frac{1}{r} \frac{\partial}{\partial r} \left(r \frac{\partial V}{\partial r} \right) + \frac{1}{r^2} \frac{\partial^2 V}{\partial \theta^2} + \frac{\partial^2 V}{\partial z^2}. \quad (\text{A.170})$$

In the spherical coordinate system, the Cartesian coordinates x , y , and z are replaced by $r = (x^2 + y^2 + z^2)^{1/2}$, $\theta = \cos^{-1}(z/r)$, and $\phi = \tan^{-1}(y/x)$. Here, r is the radial distance from the origin, θ the angle subtended between the radius vector and the z -axis, and ϕ the angle subtended between the projection of the radius vector onto the x - y plane and the x -axis. See Figure A.28. Note that r and θ in the spherical system are not the same as their counterparts in the cylindrical system. A general vector \mathbf{A} is written

$$\mathbf{A} = A_r \mathbf{e}_r + A_\theta \mathbf{e}_\theta + A_\phi \mathbf{e}_\phi, \quad (\text{A.171})$$

where $\mathbf{e}_r = \nabla r / |\nabla r|$, $\mathbf{e}_\theta = \nabla \theta / |\nabla \theta|$, and $\mathbf{e}_\phi = \nabla \phi / |\nabla \phi|$. The unit vectors \mathbf{e}_r , \mathbf{e}_θ , and \mathbf{e}_ϕ are mutually orthogonal. Hence, $A_r = \mathbf{A} \cdot \mathbf{e}_r$, et cetera. The volume element in this coordinate system is $d^3 \mathbf{r} = r^2 \sin \theta dr d\theta d\phi$. Moreover, gradient, divergence, and curl take the forms

$$\begin{aligned} \nabla V = & \frac{\partial V}{\partial r} \mathbf{e}_r + \frac{1}{r} \frac{\partial V}{\partial \theta} \mathbf{e}_\theta + \frac{1}{r \sin \theta} \frac{\partial V}{\partial \phi} \mathbf{e}_\phi, \\ \nabla \cdot \mathbf{A} = & \frac{1}{r^2} \frac{\partial}{\partial r} (r^2 A_r) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \theta} (\sin \theta A_\theta) \end{aligned} \quad (\text{A.172})$$

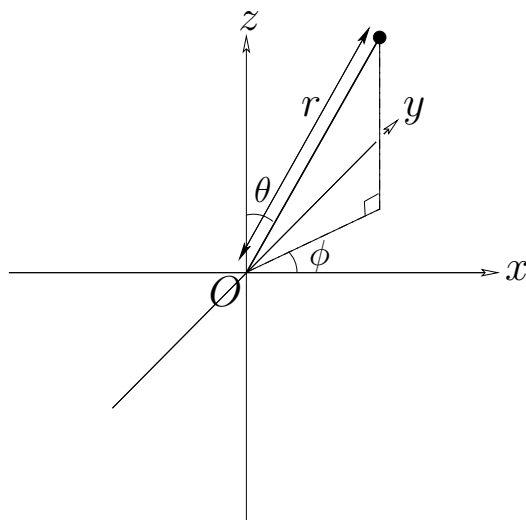


Figure A.28: Spherical coordinates.

$$+ \frac{1}{r \sin \theta} \frac{\partial A_\phi}{\partial \phi}, \quad (\text{A.173})$$

$$\begin{aligned} \nabla \times \mathbf{A} = & \left(\frac{1}{r \sin \theta} \frac{\partial}{\partial \theta} (\sin \theta A_\phi) - \frac{1}{r \sin \theta} \frac{\partial A_\theta}{\partial \phi} \right) \mathbf{e}_r \\ & + \left(\frac{1}{r \sin \theta} \frac{\partial A_r}{\partial \phi} - \frac{1}{r} \frac{\partial}{\partial r} (r A_\phi) \right) \mathbf{e}_\theta \\ & + \left(\frac{1}{r} \frac{\partial}{\partial r} (r A_\theta) - \frac{1}{r} \frac{\partial A_r}{\partial \theta} \right) \mathbf{e}_\phi, \end{aligned} \quad (\text{A.174})$$

respectively. Here, $V(\mathbf{r})$ is a general vector field, and $\mathbf{A}(\mathbf{r})$ a general scalar field. Finally, the Laplacian is written

$$\nabla^2 V = \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial V}{\partial r} \right) + \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial V}{\partial \theta} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial^2 V}{\partial \phi^2}. \quad (\text{A.175})$$

A.24 Useful Vector Identities

Notation: $\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}$ are general vectors; ϕ, ψ are general scalar fields; \mathbf{A}, \mathbf{B} are general vector fields; $(\mathbf{A} \cdot \nabla) \mathbf{B} \equiv (\mathbf{A} \cdot \nabla B_x, \mathbf{A} \cdot \nabla B_y, \mathbf{A} \cdot \nabla B_z)$ and $\nabla^2 \mathbf{A} = (\nabla^2 A_x, \nabla^2 A_y, \nabla^2 A_z)$ (but, only in Cartesian coordinates)

$$\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) = (\mathbf{a} \cdot \mathbf{c}) \mathbf{b} - (\mathbf{a} \cdot \mathbf{b}) \mathbf{c}, \quad (\text{A.176})$$

$$(\mathbf{a} \times \mathbf{b}) \times \mathbf{c} = (\mathbf{c} \cdot \mathbf{a}) \mathbf{b} - (\mathbf{c} \cdot \mathbf{b}) \mathbf{a}, \quad (\text{A.177})$$

$$(\mathbf{a} \times \mathbf{b}) \cdot (\mathbf{c} \times \mathbf{d}) = (\mathbf{a} \cdot \mathbf{c})(\mathbf{b} \cdot \mathbf{d}) - (\mathbf{a} \cdot \mathbf{d})(\mathbf{b} \cdot \mathbf{c}), \quad (\text{A.178})$$

$$(\mathbf{a} \times \mathbf{b}) \times (\mathbf{c} \times \mathbf{d}) = (\mathbf{a} \times \mathbf{b} \cdot \mathbf{d}) \mathbf{c} - (\mathbf{a} \times \mathbf{b} \cdot \mathbf{c}) \mathbf{d}, \quad (\text{A.179})$$

$$\nabla(\phi\psi) = \phi \nabla\psi + \psi \nabla\phi, \quad (\text{A.180})$$

$$\nabla(\mathbf{A} \cdot \mathbf{B}) = \mathbf{A} \times (\nabla \times \mathbf{B}) + \mathbf{B} \times (\nabla \times \mathbf{A}) + (\mathbf{A} \cdot \nabla) \mathbf{B} + (\mathbf{B} \cdot \nabla) \mathbf{A}, \quad (\text{A.181})$$

$$\nabla \cdot \nabla\phi = \nabla^2\phi, \quad (\text{A.182})$$

$$\nabla \cdot \nabla \times \mathbf{A} = 0, \quad (\text{A.183})$$

$$\nabla \cdot (\phi \mathbf{A}) = \phi \nabla \cdot \mathbf{A} + \mathbf{A} \cdot \nabla\phi, \quad (\text{A.184})$$

$$\nabla \cdot (\mathbf{A} \times \mathbf{B}) = \mathbf{B} \cdot \nabla \times \mathbf{A} - \mathbf{A} \cdot \nabla \times \mathbf{B}, \quad (\text{A.185})$$

$$\nabla \times \nabla\phi = 0, \quad (\text{A.186})$$

$$\nabla \times (\nabla \times \mathbf{A}) = \nabla(\nabla \cdot \mathbf{A}) - \nabla^2\mathbf{A}, \quad (\text{A.187})$$

$$\nabla \times (\phi \mathbf{A}) = \phi \nabla \times \mathbf{A} + \nabla\phi \times \mathbf{A}, \quad (\text{A.188})$$

$$\nabla \times (\mathbf{A} \times \mathbf{B}) = (\nabla \cdot \mathbf{B}) \mathbf{A} - (\nabla \cdot \mathbf{A}) \mathbf{B} + (\mathbf{B} \cdot \nabla) \mathbf{A} - (\mathbf{A} \cdot \nabla) \mathbf{B}. \quad (\text{A.189})$$