

Thanu Padmanabhan

Sleeping Beauties in Theoretical Physics

26 Surprising Insights

Lecture Notes in Physics

Volume 895

Founding Editors

W. Beiglböck
J. Ehlers
K. Hepp
H. Weidenmüller

Editorial Board

B.-G. Englert, Singapore, Singapore
P. Hänggi, Augsburg, Germany
W. Hillebrandt, Garching, Germany
M. Hjorth-Jensen, Oslo, Norway
R.A.L. Jones, Sheffield, UK
M. Lewenstein, Barcelona, Spain
H. von Löhneysen, Karlsruhe, Germany
M.S. Longair, Cambridge, UK
J.-M. Raimond, Paris, France
A. Rubio, Donostia, San Sebastian, Spain
M. Salmhofer, Heidelberg, Germany
S. Theisen, Potsdam, Germany
D. Vollhardt, Augsburg, Germany
J.D. Wells, Geneva, Switzerland
G. Zank, Huntsville, USA

The Lecture Notes in Physics

The series Lecture Notes in Physics (LNP), founded in 1969, reports new developments in physics research and teaching-quickly and informally, but with a high quality and the explicit aim to summarize and communicate current knowledge in an accessible way. Books published in this series are conceived as bridging material between advanced graduate textbooks and the forefront of research and to serve three purposes:

- to be a compact and modern up-to-date source of reference on a well-defined topic
- to serve as an accessible introduction to the field to postgraduate students and nonspecialist researchers from related areas
- to be a source of advanced teaching material for specialized seminars, courses and schools

Both monographs and multi-author volumes will be considered for publication. Edited volumes should, however, consist of a very limited number of contributions only. Proceedings will not be considered for LNP.

Volumes published in LNP are disseminated both in print and in electronic formats, the electronic archive being available at springerlink.com. The series content is indexed, abstracted and referenced by many abstracting and information services, bibliographic networks, subscription agencies, library networks, and consortia.

Proposals should be sent to a member of the Editorial Board, or directly to the managing editor at Springer:

Christian Caron
Springer Heidelberg
Physics Editorial Department I
Tiergartenstrasse 17
69121 Heidelberg/Germany
christian.caron@springer.com

More information about this series at <http://www.springer.com/series/5304>

Thanu Padmanabhan

Sleeping Beauties in Theoretical Physics

26 Surprising Insights



Springer

Thanu Padmanabhan
Inter-University Centre for Astronomy
and Astrophysics
Pune, Maharashtra
India

ISSN 0075-8450 ISSN 1616-6361 (electronic)
Lecture Notes in Physics
ISBN 978-3-319-13442-0 ISBN 978-3-319-13443-7 (eBook)
DOI 10.1007/978-3-319-13443-7

Library of Congress Control Number: 2015932856

Springer Cham Heidelberg New York Dordrecht London
© Springer International Publishing Switzerland 2015

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

Springer International Publishing AG Switzerland is part of Springer Science+Business Media (www.springer.com)

Preface

Theoretical physics is fun. Most of us indulge in it for the same reason a painter paints or a dancer dances — the process itself is so enjoyable! Occasionally, there are additional benefits like fame and glory and even practical uses; but most good theoretical physicists will agree that these are not the primary reasons why they are doing it. The fun in figuring out the solutions to Nature's brain teasers is a reward in itself.

The primary aim of this book is to convey this joy one feels about doing theoretical physics and share some insights in a wide variety of topics.

I recognized the need for such a book over years of teaching different aspects of theoretical physics to students and writing formal textbooks in physics. Such courses and textbooks serve a very useful purpose of training the students, but — by necessity — they cannot present the grand, unified, view of physics. Technical expertise and depth in different areas of physics comes with the price of sharp focus and detailed expositions which necessarily camouflages the broader beauty of theoretical physics. Obviously, a different kind of book — which is certainly *not* a textbook, though you might learn a lot from it — is required and I hope you find my attempt fitting the bill.

This book is a collection of 26 chapters, each devoted to highlighting some curious, fascinating and insightful aspects of a particular topic. The material ranges from a two-step (yes, exactly two steps; see Chapter 3) derivation of elliptical orbits in the inverse square law force, to regularization techniques in quantum field theory which prove that the sum of all positive integers is a negative fraction (yes; see Chapter 19). While many of the *topics* might appear to be standard, the *descriptions* are not; several professional physicists have told me that they found the discussion to be novel, many of the derivations new and the approach refreshingly different. I hope you will also find something new in this book.

Most of this book will be understandable to a bright senior undergraduate in physics who has taken basic courses in classical mechanics, quantum mechanics, special relativity and electrodynamics. I do *not* as-

sume previous acquaintance with quantum field theory or general relativity (though some of the chapters deal with these topics). *You can dip in anywhere you please in this book and start reading!* The chapters are reasonably modular (except for a few obvious ones which come in pairs). You will find the highlights of each of the chapters described just after the table of contents which will help you to decide how you want to proceed. Further, instead of subsections, I have sprinkled marginal comments throughout the book which will alert you as to what is being talked about in the corresponding paragraph; this makes the book even more modular to use! You will find a list of references right at the end which could guide you for further reading, although virtually every topic discussed here can be pursued further by simple web-based searches. Partly for this reason, I have kept the references rather minimal and I apologize to anyone whose contribution might have been overlooked.

Many people have contributed, in different ways, to the making of this book. Angela Lahee of Springer initiated this project and helped me through its completion, displaying considerable initiative. Several of the chapters overlap in their intellectual content with a series of articles I wrote in the journal *Resonance* during 2008-2009 even though they have all undergone significant amount of re-writing, re-grouping and inclusion of additional material and topics. I thank the Indian Academy of Sciences for granting permission to Springer for the reuse of the material in these articles in this book. Many of my colleagues went through the previous drafts of the book and offered comments. Special thanks are due to Hamsa Padmanabhan and Aseem Paranjape for detailed comments and corrections in several chapters. I thank the following colleagues (listed in alphabetical order) for comments on different chapters in the earlier drafts: Jasjeet Bagla, Prasanta Bera, Pallavi Bhat, Sumanta Chakraborty, George Ellis, Bhooshan Gadre, Peter Goldreich, Neeraj Gupta, Nissim Kanekar, Vikram Khair, Dawood Kothawala, Kinjalk Lochan, Malcolm Longair, Abhilash Mishra, Dipanjan Mukherjee, Suvodip Mukherjee, Krishamohan Parattu, Tirthankar Roy Choudhury, Kanak Saha, Sudipta Sarkar, S. Shankaranarayanan, Suprit Singh, T.P. Singh, Kandaswamy Subramanian, Durgesh Tripathi.

This book would not have been possible without the dedicated support from Vasanthi Padmanabhan, who not only did the entire LaTeXing and formatting but also produced most of the figures. I thank her for her help. It is a pleasure to acknowledge the library and other facilities available at IUCAA, which were useful in this task.

Pune,
September 2014

Thanu Padmanabhan

Contents

Chapter Highlights	IX
Notations and Conventions	XVI
1 The Grand Cube of Theoretical Physics	1
2 The Emergence of Classical Physics	7
3 Orbits of Planets <i>are</i> Circles!	25
4 The Importance of being Inverse-square	43
5 Potential surprises in Newtonian Gravity	57
6 Lagrange and his Points	65
7 Getting the most of it!	73
8 Surprises in Fluid Flows	89
9 Isochronous Curiosities: Classical and Quantum	99
10 Logarithms of Nature	109
11 Curved Spacetime for pedestrians	117
12 Black hole is a Hot Topic	127
13 Thomas and his Precession	135
14 When Thomas met Foucault	143
15 The One-body Problem	153

16 The Straight and Narrow Path of Waves 167

17 If Quantum Mechanics is the Paraxial Optics, then 175

18 Make it Complex to Simplify 191

19 Nothing matters a lot 205

20 Radiation: Caterpillar becomes Butterfly 219

21 Photon: Wave and/or Particle 231

22 Angular Momentum without Rotation 241

23 Ubiquitous Random Walk 247

24 More on Random Walks: Circuits and a Tired Drunkard 259

25 Gravitational Instability of the Isothermal Sphere 269

26 Gravity bends electric field lines 279

References 293

Index 299

Chapter Highlights

1. The Grand Cube of Theoretical Physics

The ‘big picture’ of theoretical physics can be nicely summarized in terms of a unit cube made of the fundamental constants G, \hbar, c^{-1} representing the three axes. The vertices and linkages of this cube — which we will explore in different chapters of this book — allow you to appreciate different phenomena and their inter-relationships. This chapter introduces the Cube of Theoretical Physics and relates it to the rest of the book.

2. The Emergence of Classical Physics

Quantum physics works with probability amplitudes while classical physics assumes deterministic evolution for the dynamical variables. For example, in non-relativistic quantum mechanics, you will solve the Schrodinger equation in a potential to obtain the wave function $\psi(t, q)$, while the same problem — when solved classically — will lead to a trajectory $q(t)$. How does a deterministic trajectory arise from the foggy world of quantum uncertainty? We will explore several aspects of this correspondence in this chapter, some of which are nontrivial. You will discover the *real* meaning of the Hamilton-Jacobi equation (without the usual canonical transformations, generating functions and other mumbo-jumbo) and understand why the Hamilton-Jacobi equation told us $p_a = \partial_a S = (-\partial_t S, \nabla S) = (E, \mathbf{p})$ even before the days of four vectors and special relativity. We will also address the question of why the Lagrangian is equal to kinetic energy minus potential energy (or is it, really?) and why there are *only* two classical fields, electromagnetism and gravity. In fact, you will see that classical physics makes better sense as a limit of quantum physics!

3. **Orbits of Planets *are* Circles!**

The orbits of planets, or any other body moving under an inverse square law force, can be understood in a simple manner using the idea of the velocity space. Surprisingly, a particle moving in an ellipse, parabola or a hyperbola in real space moves in a circle in the velocity space. This approach allows you to solve the Kepler problem in just two steps! We will also explore the peculiar symmetry of the Lagrangian that leads to the conservation of the Runge-Lenz vector and the geometrical insights that it provides. Proceeding to the relativistic versions of Kepler/Coulomb problem you will discover why the forces *must be* velocity dependent in a relativistic theory and describe a new feature in the special relativistic Coulomb problem, viz. the existence of orbits spiraling to the center.

4. **The Importance of being Inverse-square**

This chapter continues the exploration started in the previous one. The Coulomb problem, which corresponds to motion in a potential that varies as r^{-1} , has a peculiar symmetry which leads to a phenomenon known as ‘accidental’ degeneracy. This feature exists both in the classical and quantum domains and allows some interesting, alternative ways to understand, e.g., the hydrogen atom spectrum. We will see how one can find the energy levels of the hydrogen atom without solving the Schrodinger equation and how to map the 3D Coulomb problem to a 4D harmonic oscillator problem. The $(1/r)$ nature of the potential also introduces several peculiarities in the *scattering* problem and we will investigate the questions: (i) How come quantum Coulomb scattering leads *exactly* to the Rutherford formula? What happened to the \hbar ? (ii) How come the Born approximation gives the exact result for the Coulomb potential? What do the ‘unBorn’ terms contribute?!

5. **Potential surprises in Newtonian Gravity**

How unique is the distribution of matter which will produce a given Newtonian gravitational field in a region of space? For example, can a non-spherical distribution of matter produce a strictly inverse square force outside the source? Can a non-planar distribution of matter produce a strictly constant gravitational force in some region? We discuss the rather surprising answers to these questions in this chapter. It turns out that the relation between the density distribution and the gravitational force is far from what one would have naively imagined from the textbook examples.

6. **Lagrange and his Points**

A solution to the 3-body problem in gravity, due to Lagrange, has several remarkable features. In particular, it describes a situation in

which a particle, located at the *maxima* of a potential, remains stable against small perturbations. We will learn a simple way of obtaining this equilateral solution to the three body problem and understanding its stability.

7. Getting the most of it!

Extremum principles play a central role in theoretical physics in many guises. We will discuss, in this chapter, some curious features associated with a few unusual variational problems. We start with a simple way to solve the standard brachistochrone problem and address the question: How come the cycloid solves all the chron-ic problems? (Or does it, really?). We then consider the brachistochrone problem in a real, $(1/r^2)$, gravitational field and describe a new feature which arises: viz. the existence of a forbidden zone in space not accessible to brachistochrone curves! We will also determine the shape of a planet that exerts the maximum possible gravitational force at a point on its surface — a shape which does not even have a name! Finally, we take up the formation of the rainbows with special emphasis on the question: Where do you look for the tertiary (3rd order) rainbow?

8. Surprises in Fluid Flows

The idealized flow of a fluid around a body is a classic text book problem in fluid mechanics. Interestingly enough, it leads to some curious twists and conceptual conundrums. In particular, it leads a surprising divergence which needs to be regularized even in the text book case of fluid flow past a sphere!

9. Isochronous Curiosities: Classical and Quantum

The oscillatory motion of a particle in a one dimensional potential belongs to a class of exactly solvable problems in classical mechanics. This chapter examines some lesser known aspects of this problem in classical and quantum mechanics. It turns out that both $V(x) = ax^2$ and $V(x) = ax^2 + bx^{-2}$ have (1) periods of oscillation which are independent of amplitude in classical physics and (2) equally spaced energy levels in quantum theory. We will explore several features of this curious correspondence. We will also discuss the question of determining the potential from the period of oscillation (in classical physics) or from the energy levels (in quantum physics) which are closely related and clarify several puzzling features related to this issue.

10. Logarithms of Nature

Scaling arguments and dimensional analysis are powerful tools in physics which help you to solve several interesting problems. And when the scaling arguments fail, as in the examples discussed in this

chapter, we are led to a more fascinating situation. A simple example in electrostatics leads to infinities in the Poisson equation and we get a finite E from an infinite ϕ ! I also describe the quantum energy levels in the delta function potentials and show how QFT helps you to understand QM better!

11. **Curved Spacetime for pedestrians**

The spacetime around a spherical body plays a key role in general relativity and is used in the crucial tests of Einstein's theory of gravity. This spacetime geometry is usually obtained by solving Einstein's equations. I will show how this metric can be obtained by a simple — but strange — trick. Along the way, you will also learn a three-step proof as to why gravity must be geometry, the reason why the Lagrangian for a particle in a Newtonian gravitational field is kinetic energy minus potential energy and how to obtain the orbit equation in GR, just from the principle of equivalence.

12. **Black hole is a Hot Topic**

A fascinating result in black hole physics is that they are not really black! They glow as though they have a surface temperature which arises due to purely quantum effects. I will provide a simple derivation of this hot result based on the interpretation of a plane wave by different observers.

13. **Thomas and his Precession**

Thomas precession is a curious effect in special relativity which is purely kinematical in origin. But it illustrates some important features of the Lorentz transformation and possesses a beautiful geometric interpretation. We will explore the physical reason for Thomas precession and its geometrical meaning in this chapter and in the next.

14. **When Thomas met Foucault**

The Foucault pendulum is an elegant device that demonstrates the rotation of the Earth. We describe a paradox related to the Foucault pendulum and provide a geometrical approach to determine the rotation of the plane of the pendulum. By introducing a natural metric in the velocity space we obtain an interesting geometrical relationship between the dynamics of the Foucault pendulum and the Thomas precession discussed in the previous chapter. This approach helps us to understand both phenomena better.

15. **The One-body Problem**

You might have thought that the one-body problem in physics is trivial. Far from it! One can look at the free particle in an inertial or a non-inertial frame, relativistically or non-relativistically, in flat or

in curved spacetime, classically or quantum mechanically. All these bring in curious correspondences in which the more exact theory provides valuable insights about the approximate description. I start with the surprising — and not widely appreciated — result that you really can't get a sensible free-particle Lagrangian in non-relativistic mechanics while you can do it in relativistic mechanics. In a similar vein, the solution to the Klein-Gordon equation transforms as a scalar under coordinate transformations, while the solution to the Schrodinger equation does not. These conundrums show that classical mechanics makes more sense as a limiting case of special relativity and the non-relativistic Schrodinger equation is simpler to understand as a limiting case of the relativistic Klein-Gordon equation!

16. **The Straight and Narrow Path of Waves**

Discovering unexpected connections between completely different phenomena is always a delight in physics. In this chapter and the next, we will look at one such connection between two unlikely phenomena: propagation of light and the path integral approach to quantum *field* theory! This chapter introduces the notion of paraxial optics in which we throw away half the solutions and still get useful results! I also describe the role of optical systems and how the humble lens acts as an analog device that performs Fourier transforms. In passing, you will also learn how Faraday's law leads to diffraction of light.

17. **If Quantum Mechanics is the Paraxial Optics, then**

The *quantum mechanical* amplitude for a particle to propagate from event to event in spacetime shows some nice similarities with the corresponding propagator for the electromagnetic wave amplitude discussed in the previous chapter. This raises the question: If quantum mechanics is paraxial optics, what is the exact theory you get when you go beyond the paraxial approximation? In the path integral approach to quantum mechanics you purposely avoid summing over *all* the paths while in the path integral approach for a *relativistic* particle you *are forced to* sum over all paths. This fact, along with the paraxial optics analogy, provides an interesting insight into the transition from quantum field theory to quantum mechanics and vice versa! I also describe why combining the principles of relativity and quantum theory *demand*s a description in terms of fields.

18. **Make it Complex to Simplify**

Some of the curious effects in quantum theory and statistical mechanics can be interpreted by analytically continuing the time coordinate to purely imaginary values. We explore some of these issues in this chapter. In quantum mechanics, this allows us to determine the properties of ground state from an approximate evaluation of path integrals. In statistical mechanics this leads to an unexpected connection

between periodicity in imaginary time and temperature. The power of this approach can be appreciated by the fact that one can derive the black hole temperature in just a couple of steps using this procedure. Another application of the imaginary time method is to understand phenomena like the Schwinger effect which describes the popping out of particles from the vacuum. Finally, I describe a non-perturbative result in quantum mechanics, called the over-the-barrier-reflection, which is easier to understand using complex paths.

19. Nothing matters a lot

The vacuum state of the electromagnetic field is far from trivial. Amongst other things, it can exert forces that are measurable in the lab. This curious phenomenon, known as the *Casimir effect*, is still not completely understood. I describe how the probability distribution for the existence of electromagnetic fields in the vacuum can be understood, just from the knowledge of the quantum mechanics of the harmonic oscillator. This chapter also introduces you to the tricks of the trade in quantum field theory, which are essential to get finite answers from divergent expressions - like to prove that the sum of all positive integers is a negative fraction!

20. Radiation: Caterpillar becomes Butterfly

The fact that an accelerated charge will radiate energy is considered an elementary textbook result in electromagnetism. Nevertheless, this process of radiation (and its reaction on the charged particle) raises several conundrums about which technical papers are written even today. In this chapter, we will try to understand how the caterpillar ($1/r^2$, radial field) becomes a butterfly ($1/r$, transverse field) in a simple, yet completely rigorous, manner without the Lienard-Wiechert potentials or other red-herrings. I will also discuss some misconceptions about the validity of $\nabla \cdot \mathbf{E} = 4\pi\rho$ for radiative fields with retardation effects.

21. Photon: Wave and/or Particle

The interaction of charged particles with blackbody radiation is of considerable practical and theoretical importance. Practically, it occurs in several astrophysical scenarios. Theoretically, it illustrates nicely the fact that one can think of the radiation either as a bunch of photons or as electromagnetic waves and still obtain the same results. We shall highlight some non-trivial aspects of this correspondence in this chapter. In particular we will see how the blackbody radiation leads a double life of being either photons or waves and how the radiative transfer between charged particles and black body radiation can be derived just from a Taylor series expansion (and a little trick)! Finally, I will describe the role of radiation reaction force on charged particles to understand some of these results.

22. Angular Momentum without Rotation

Not only mechanical systems, but also electromagnetic fields carry energy and momentum. What is not immediately apparent is that certain static electromagnetic configurations (with no rotation in sight) can also have angular momentum. This leads to surprises when this angular momentum is transferred to the more tangible rotational motion of charged particles coupled to the electromagnetic fields. A simple example described in this chapter illustrates, among other things, how an observable effect arises from the unobservable vector potential and why we can be cavalier about gauge invariance in some circumstances.

23. Ubiquitous Random Walk

What is common to the spread of mosquitoes, sound waves and the flow of money? They all can be modeled in terms of random walks! Few processes in nature are as ubiquitous as the random walk which combines extraordinary simplicity of concept with considerable complexity in the final result. In this and the next chapter, we shall examine several features of this remarkable phenomenon. In particular, I will describe the random walk in the velocity space for a system of gravitating particles. The diffusion in velocity space can't go on and on — unlike that in real space — which leads to another interesting effect known as dynamical friction — first derived by Landau in an elegant manner.

24. More on Random Walks: Circuits and a Tired Drunkard

We continue our exploration of random walks in this chapter with some more curious results. We discuss the dimension dependence of some of the features of the random walk (e.g., why a drunken man will eventually come home but a drunken bird may not!), describe a curious connection between the random walk and electrical networks (which includes some problems you can't solve by being clever) and finally discuss some remarkable features of the random walk with decreasing step-length, which is still not completely understood and leads to Cantor sets, singularities and the golden ratio — in places where you don't expect to see them.

25. Gravitational Instability of the Isothermal Sphere

The statistical mechanics of a system of particles interacting through gravity leads to several counter-intuitive features. We explore one of them, called Antonov instability, in this chapter. I describe why the thermodynamics of gravitating systems is non-trivial and how to obtain the mean-field description of such a system. This leads to a self-gravitating distribution of mass called the isothermal sphere which exhibits curious features both from the mathematical and physical

points of view. I provide a simple way of understanding the stability of this system, which is of astrophysical significance.

26. Gravity bends electric field lines

Field lines of a point charge are like radially outgoing light rays from a source. You know that the path of light is bent by gravity; do electric field lines also bend in a gravitational field? Indeed they do, and — in the simplest context of a constant gravitational field — both are bent in the same way. Moreover, both form arcs of circles! The Coulomb potential in a weak gravitational field can be expressed in a form which has unexpected elegance. The analysis leads to a fresh insight about electromagnetic radiation as arising from the weight of electrostatic energy in the rest frame of the charged particle, and also allows you to obtain Dirac's formula for the radiation reaction, in three simple steps.

Notations and Conventions

Most of the notations used in the book are fairly standard. You may want to take note of the following:

1. I use the Gaussian system of units to describe electromagnetic phenomena; however, conversion to SI units is completely straightforward in all the relevant chapters.
2. In chapters involving relativity, the Latin letters a, b, \dots range over the spacetime indices 0, 1, 2, 3, while the Greek indices α, β, \dots range over the spatial coordinates 1, 2, 3 with the notation $\partial_i = (\partial/\partial x^i)$ for coordinate derivatives. When the discussion does not involve relativistic physics, this distinction between Latin and Greek subscripts is not maintained. The signature for the spacetime is $(-, +, +, +)$ with $\eta_{ij} = \text{dia}(-1, 1, 1, 1) = \eta^{ij}$. Units with $c = 1$ are used most of the time though c is re-introduced when required.
3. All through the book (and not only in chapters dealing with relativity) I use the summation convention according to which any index repeated in an algebraic expression is summed over its range of values.
4. In topics dealing with quantum mechanics, I often use units with $\hbar = 1$, re-introducing it into the equations only when relevant.
5. In the equations, you will sometimes find the use of the symbol \equiv . This indicates that the equation defines a new variable or notation.

The Grand Cube of Theoretical Physics

1

The key purpose of this book is to let you *enjoy* theoretical physics, appreciate the beautiful overall structure and see how everything hangs together. To do this, it is helpful to have a map which will allow you to navigate the landscape of theoretical physics. The Cube of Theoretical Physics (CTP) — which I will now introduce — is a good way to begin.

*The landscape of
Theoretical Physics*

The fundamental principles of physics emphasize the role of three constants: G (Newton's gravitational constant), c (the speed of light) and \hbar (the Planck constant). By a suitable choice of units, we can set the numerical value of each of these three to unity and the broad structure of physical theories can be described using a 3-dimensional space in which each of the Cartesian coordinates (see Fig. 1.1) is taken to be one of the above mentioned fundamental constants. (I have used this diagram during my lectures in the mid-eighties. A somewhat similar diagram with a tetrahedron rather than a cube appears in Ref. [1]. It is very likely that many others have thought of such a description but the only published reference I know is Ref. [2] of which I am a co-author.) It is convenient to use $(1/c)$ rather than c in such a description. The entire space of physical theories will be confined within the unit cube so formed.

An examination of this diagram reveals several interesting features. The origin $G = 0$, $\hbar = 0$, $c^{-1} = 0$ represents an idealized non-relativistic (point) mechanics (NRM) with which your physics course begins. Starting from this and traveling along different directions on the CTP, we can have a glimpse of what nature has in store.

Moving along the G -axis will lead to non-relativistic, classical, Newtonian gravity (NG) which is probably the least disturbing journey one can undertake on the CTP. In fact, your classical mechanics course starting at the origin will certainly include topics like the Kepler problem which uses Newtonian gravity. So the vertical axis between NRM and NG in Fig. 1.1 is usually treated together in your first semester course. To be honest, this completely hides the true nature of gravity but then, as we will see repeatedly, physicists love useful approximations.

*Newtonian gravity +
Classical Mechanics*

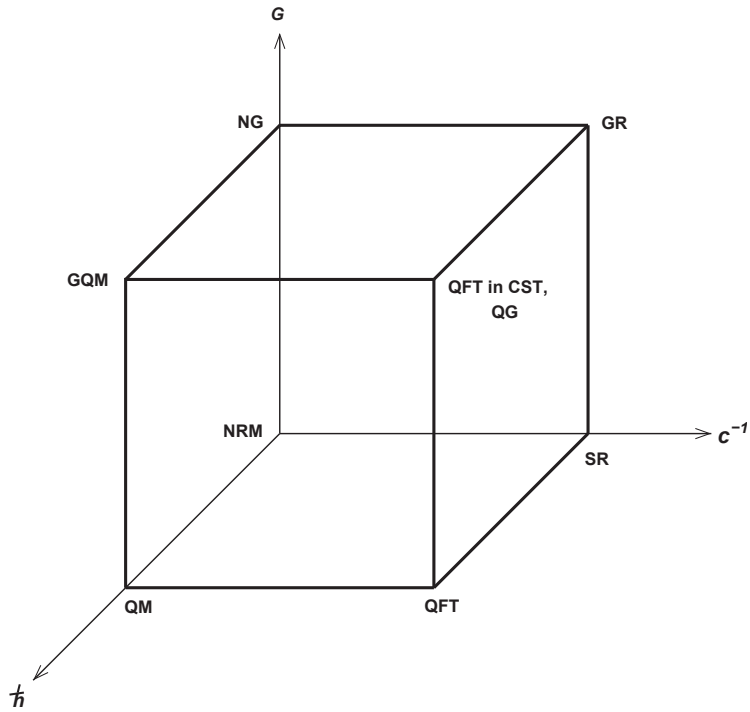


Fig. 1.1: The landscape of theoretical physics can be concisely described by a cube — The Cube of Theoretical Physics — whose axes represents the three fundamental constants G, \hbar and c^{-1} . The vertices and linkages describe different structural properties of the physical theories. See text for detailed description.

Space + Time = Spacetime

Moving along the speed of light axis to $c^{-1} = 1$, (keeping $G = 0, \hbar = 0$), will get you to special relativistic (SR) mechanics. Instead of space and time being treated as separate entities, we now view them as parts of a spacetime continuum. Time is no longer absolute and two clocks moving with respect to each other run at different rates.

Real world is quantum!

Traveling along the \hbar –axis will lead to non-relativistic quantum mechanics (QM) and, as they say, if you are not shocked on the first exposure to quantum mechanics, you haven’t grasped it! It turns the deterministic classical world on its head and introduces probabilistic concepts, wave-functions, wave-particle duality and all the rest. You slowly get used to it.

Gravity is just spacetime geometry!

These three vertices provide more accurate descriptions of nature than the region near the origin but it gets better if we keep a *pair* of constants to be non-zero. The vertex $c^{-1} = 1, G = 1, \hbar = 0$ represents classical general relativity (GR) which combines the principles of special relativity and gravity. This takes you from flat spacetime to curved spacetime and

tells you that gravity is actually a manifestation of curved geometry. In SR, you learnt that clocks in relative motion will run at different rates with respect to each other; now you learn that even clocks at rest with respect to each other can run at different rates, if they are located at different gravitational potentials. Gravity affects the flow of time!

Similarly, $\hbar = 1$, $c^{-1} = 1$, $G = 0$ leads to flat spacetime quantum field theory (QFT), which combines the principles of special relativity and quantum theory. Particles lose their eternal existence and can now pop in and out of the vacuum. Further, it is mandatory that every particle must have an antiparticle and that interactions are mediated by the exchange of special kinds of particles. In fact, you need a totally new kind of language to understand these high energy phenomena.

Particles must have antiparticles!

The vertex at which all the three constants are unity, $c^{-1} = 1$, $G = 1$, $\hbar = 1$, should represent the domain of quantum gravity but — more importantly for our purpose — *it also represents the study of quantum field theory in curved space-time* (QFT in CST), like, for example, the study of radiation from black holes. A description of the thermal features of black holes (Chapter 12) requires all these three constants to be non-zero. While quantum gravity still remains a distant dream, we do have a fair amount of understanding of quantum field theory in curved spacetime and, in this sense, this vertex (QFT in CST) can be considered to be within our grasp.

*The final frontier:
Gravity + Quantum*

While most of the above limiting forms of physical theories have attracted a reasonable amount of attention and made it into textbooks, the Fig. 1.1 brings out one limiting case which probably has *not* been explored in comparable detail [2]. This is the “ignored” vertex with $c^{-1} = 0$, $G = 1$, $\hbar = 1$, which corresponds to exploring the nature of gravity in a *quantum mechanical* context (GQM). Some of the discussion in a later chapter, Chapter 15, will be devoted to the exploration of this vertex. More generally, we will deal with the issue of projecting theories to the $G\hbar$ plane by taking the $c \rightarrow \infty$ limit in different contexts.

The ignored vertex of theoretical physics!

The chapters of this book will take you through a tour of the CTP. There are a few chapters which will linger on a particular vertex just to explore some curious features there. And then there are other chapters which tell you what happens on the links when the effects of two vertices are incorporated or describe the curious limiting behaviour as we go from different vertices towards the origin.

Sneak Preview of the book

As you can easily imagine most of the topics require inputs from more than one vertex and hence sit on the linkages. For example, Chapter 3 (where you learn that planets actually move in circular orbits) provides a two-step solution to Kepler problem and demystifies several aspects of it. This clearly sits on the link between NRM and NG. The closely related Chapter 4, involving NRM and QM, studies quantum mechanical aspects of the Kepler/Coulomb problem and shows you that the hydrogen atom is essentially a harmonic oscillator in disguise. Chapter 5 (where you learn that planets of weird, non-spherical shape, can exert a *strictly* $1/r^2$ force

The hidden charms of Newtonian gravity

outside — a result which many physicists feel is impossible) and Chapter 6 (which tells you how motion can be perfectly stable around *maxima* of the potential and how Nature exploits this result) use both NRM and NG and belong to that link. Chapter 25 also investigates the NG-NRM link in the context of thermodynamics of gravitating systems, which is nothing like the thermodynamics of the usual gaseous systems you would have learnt in standard courses.

Chapter 9 studies a special class of potentials in which the classical period of oscillation is independent of the amplitude and explores its quantum analogues drawing from both QM and NRM. Chapters 10 and 17 possibly belong to the QM-QFT link. Chapter 10 illustrates the ideas of the renormalization group in an elementary example from quantum mechanics. Chapter 17 describes the relation between quantum mechanics and optics and shows how one can understand the transition from QM to QFT exploiting a simple optics analogy.

There are two chapters dealing with the manner in which approximate descriptions emerge from more exact descriptions. Chapter 2 tells you how trajectories of particles arise in Newtonian, special relativistic, and even general relativistic physics from corresponding quantum descriptions. This belongs to at least three links, GQM-QM-NRM, QFT-SR-NRM and QFT in CST-GR-NG. Chapter 15 explains why you need special relativity if you have to understand non-relativistic mechanics properly!

Other chapters can be mostly confined to a single vertex of CTP. Chapter 7 (which is a potpourri of extremum problems including the brachistochrone in an inverse square force field, the strange shape of a planet that can exert the maximum possible force at a point on its surface, and why it is so hard to see the tertiary rainbow), Chapter 8 (which tells you how strange conundrums can arise in the simplest of the fluid flows), Chapter 23 (introducing the concepts of dynamical friction and velocity relaxation in stellar systems) and Chapter 24 (where you explore unexpected features of random walks like their relation to electric circuits and how a drunkard who is getting progressively tired can lead you to a Cantor set) are probably closest to your standard classical mechanics course, and they live at the NRM vertex.

There are several chapters which deal with the SR vertex. Chapter 13 describes a phenomenon called Thomas precession which is counterintuitive but has a lovely geometrical interpretation. Surprisingly, the mathematics is essentially the same as that of the Foucault pendulum — a connection which you might not have suspected *a priori*. This is described in Chapter 14 which probably falls somewhere along the SR-GR link. While we are not using curved *spacetime*, some notions of curved geometry (in the velocity space!) find application here. I will put Chapter 22 (which describes a perfectly static electromagnetic field filled to the brim with angular momentum), Chapter 20 (where we learn how to get the ex-

*The enrichment
from QM*

*The non-trivial
limits*

*A potpourri of
curious physics*

*Special relativity
and electromag-
netism*

act electromagnetic fields of an arbitrarily moving charged particle without differentiating the Lienard-Wiechert or anybody else's potential), and Chapter 16 (which describes optics in a manner that will be useful later on to explore the connection between quantum mechanics and quantum field theory) also at the SR vertex. This is because anything electromagnetic properly belongs to the domain of special relativity. (There are, alas, textbooks which will begin to teach you electrodynamics without special relativity and bring it in after a dozen chapters; if you learnt electrodynamics from one of them, may be you need a remedial course!)

I have included two chapters dealing with the GR vertex, viz., Chapters 11 and 26. Chapter 11 shows you how to get the curved spacetime around a spherical body by a cute trick — which works for reasons nobody really understands. Chapter 26 discusses how gravity bends the electric field lines of a charged particle and shows you that, in the simplest context, this bending of electric field lines is exactly the same as the bending of light rays by gravity!

General relativity

Two Chapters (18 and 19) are explorations in quantum field theory. One deals with the fascinating manifestation of vacuum fluctuations known as the Casimir effect, which describes the force of attraction between two conducting plates kept in the vacuum; along the way, you learn that the sum of all positive integers is actually a negative fraction, viz. $(-1/12)$ (incredible, but true!). The other deals with the production of particles from the vacuum and shows how it can be thought of as due to complex trajectories of virtual particles. Chapter 21 explores the interaction of charged particles with radiation when the latter is treated either as fluctuating electromagnetic fields or as a bunch of photons, and elucidates the wave-particle duality as applied to the photon in a very practical context.

Quantum field theory

The exploration of black hole thermodynamics (Chapter 12), possibly the only concrete result we have in combining the principles of gravity and quantum mechanics, belongs to the diagonally opposite vertex to the origin (viz. QFT in CST). I provide an accessible, simple, yet rigorous, derivation this result.

When QM met GR

The tour around CTP also highlights the following amusing fact: It is incredible how generations of theoretical physicists are trained, starting from a model of the world which is known to be completely wrong! Semester after semester you correct and relearn the wrong things you have learnt before. After a course in classical mechanics, you will be told that there is something called special relativity and the Newton's laws are wrong. You will then learn that when gravity is included, special relativity is no good and you need to redo everything in curved spacetime to include gravitational physics, because Newton got not only his equations of motion wrong but also his law of gravitation. While you are grappling with all these some other professor would have told you that even the entire fabric of physics you have been taught in previous semesters is incorrect

The way we learn physics!

and that the world is (something loosely described as) quantum mechanical. There is no deterministic evolution and everything has to be done in a probabilistic manner. You learn that all the physics you have learnt (except thermodynamics, but we will not get into that) needs to be quantized — which might take up couple of more semesters. If you still persist with physics, you will learn how to put together special relativity and quantum mechanics in the form of quantum field theory and maybe even learn how to do field theory in a curved background, thereby bringing together gravity and quantum mechanics in a rough sort of way. Clearly, education in advanced physics is a progressive attempt to correct the wrong things taught to you earlier!

Some physicists will protest and say, “Well, you see, it is not really wrong physics we teach; it is all valid in some approximate sense. Anyway, a student cannot understand advanced concepts all at one go. It has to be given in small doses, one step at a time”. There *is* lot of practical truth in this claim but one cannot but notice that no mathematician is ever taught anything wrong (or approximate) — but we physicists learn to live with approximations and idealizations which get corrected progressively. This is the price we pay to be able to relate to real Nature out there (which pure mathematics is not overly concerned with!). Hopefully this book will also help you to appreciate the broader structure of theoretical physics and how approximations are embedded in more exact descriptions.

Read on!
And have fun!

The Emergence of Classical Physics

2

Quantum physics is nothing like classical physics and it is probably not an exaggeration to say that we just get used to quantum physics — without really understanding it — as we learn more about it! There are several conceptual and technical problems involved in taking the classical limit of a quantum mechanical description. We will not worry too much about the conceptual issues — interesting though they are — but will instead concentrate on one technical issue in this chapter.

*Quantum World:
amplitudes, probabilities and
uncertainties*

The central quantity in quantum physics is the probability amplitude for something to happen, described by a complex number Ψ . In the simplest case of non-relativistic quantum mechanics, this could be the wavefunction $\psi(t, q)$ for a particle such that $|\psi(t, q)|^2$ gives the probability to find this particle at a position q at time t . The same kind of idea works even in more general contexts. For example, one can study the quantum version of electrodynamics in terms of a similar amplitude $\Psi(\mathbf{E}(\mathbf{x}), t)$ such that $|\Psi(\mathbf{E}(\mathbf{x}), t)|^2$ gives the probability that an electric field $\mathbf{E}(\mathbf{x})$ exists in space at time t . (We will say more about this in Chapter 19.) In all these cases, the amplitude satisfies a *linear* equation allowing the superposition of solutions of the equation. In the case of non-relativistic quantum mechanics, this is just the Schrödinger equation; in more complex cases the equation can be more complicated but is *always* linear in the amplitude.

Classically, on the other hand, we describe the same system by a deterministic evolution. In non-relativistic mechanics, our aim is to find the trajectory $q(t)$ of a particle, by solving, say, Newton's law of motion; in classical electrodynamics, we determine the evolution of the electric field at all times by finding $\mathbf{E}(t, \mathbf{x})$ as a solution to Maxwell's equations. No probabilities, no probability amplitudes! How do we get here from there?

*Classical World:
trajectories, deterministic evolution*

The answer is fascinating and, in fact, validates several techniques used in classical physics that appear contrived or mysterious within the classical context. Let me first explain qualitatively how this comes about.

The key idea is to write the quantum amplitude — which is a complex number — in the form:

$$\Psi = R \exp\left(\frac{iS}{\hbar}\right), \quad (2.1)$$

which is just the standard representation of a complex number in terms of an amplitude and a phase — with the crucial new input being the way we have introduced \hbar (we will say more about it soon). This way of representing the quantum amplitude gives us a clue as to how the classical physics might arise. If we substitute this expression into the equation satisfied by the amplitude (which is just the Schrödinger equation in the case of non-relativistic quantum mechanics, or a more complicated one in other contexts) and equate the real and imaginary parts, we will obtain two equations for R and S — which, of course, are completely equivalent to the original equation satisfied by the amplitude Ψ . We now assume that the phase S is analytic in \hbar and has a Taylor series expansion:

$$S = S_0 + \hbar S_1 + \hbar^2 S_2 + \dots \quad (2.2)$$

which means that, at the lowest order, the phase of Ψ in Eq. (2.1) is given by (S_0/\hbar) and is *non-analytic* in \hbar .

Incredibly enough, we can solve the relevant equation (which is the Schrödinger equation in non-relativistic quantum mechanics) consistently, order by order in \hbar and — in particular — determine S_0 , which is independent of \hbar . The fact that S_0 satisfies an equation that is independent of \hbar not only in non-relativistic quantum mechanics *but in all physical theories known to us* is quite non-trivial. It tells you something deep about the laws of nature.

When we solve these equations, we will introduce some additional constants (analogous to integration constants) in the solution. Let us denote one such constant by λ and the corresponding lowest order phase by (S_λ/\hbar) , which depends on λ . (We have dropped the subscript 0 in S_0 for notational simplicity and written $S_0(\lambda) = S_\lambda$.) Then the probability amplitude will depend on λ and one could write $\Psi_\lambda = R \exp(iS_\lambda/\hbar)$ for the particular solution correct to the lowest order. (Strictly speaking, we should use the notation R_λ rather than R , but it will turn out that R plays only a minor role in what follows; so we will not bother about it.) But since the original equation satisfied by Ψ is linear in Ψ , one can superpose solutions with different λ to find a general solution. When we add the solutions with different λ , we are adding waves with different phases (S_λ/\hbar) . (Again, strictly speaking, the amplitudes R_λ are also dependent on λ but this dependence is irrelevant for the interference condition at the leading order.) In the limit of $\hbar \rightarrow 0$, the phases will oscillate rapidly and waves with different values of λ will cancel each other out in general. We will get a non-zero result only if the phase does not change significantly

The phase of the quantum amplitude

Classical world from the constructive interference of quantum waves!

for small changes in λ . *This condition for stationary phase translates to $(\partial S/\partial \lambda) = 0$ which selects out a classical evolutionary history!* In non-relativistic quantum mechanics, for example, $S = S(t, q, \lambda)$ and the condition $\partial S/\partial \lambda = 0$ will lead to a trajectory $q = q(t, \lambda)$. Thus, the classical trajectories arise from the condition for the stationarity of the phase of the quantum wavefunction. All we need to check, of course, is that this does give the expected classical trajectory.

As I said before, *everything we know in classical physics arises from the corresponding quantum description by the above mechanism.* We will first try this out in the context of non-relativistic quantum mechanics and explore some nuances, before describing more general cases.

In the context of a non-relativistic particle influenced by a potential $V(q)$, the amplitude $\psi(t, q)$ satisfies the time-dependent Schrödinger equation

*Quantum to
Classical: non-
relativistic particle*

$$i\hbar\dot{\psi} = -\frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial q^2} + V(q)\psi, \quad (2.3)$$

where overdot denotes derivative with respect to time. You also know that, classically, the same particle is described by a Hamiltonian $H(p, q)$, and an equation of motion:

$$H(p, q) = \frac{p^2}{2m} + V(q); \quad m\ddot{q} = -V'(q). \quad (2.4)$$

In fact, you learn the wrong theory (classical physics) first and then ‘quantize’ it to get a better description — in this case, through the Schrödinger equation in Eq. (2.3) obtained from $H(p, q)$. But let us forget this historical fact and assume that you are just given the more accurate theory, in the form of Eq. (2.3). You know that the classical behaviour — trajectories and all — has to emerge from this equation in the limit of $\hbar \rightarrow 0$. How do we go about taking this limit?

It is worth thinking about this issue a little bit more before jumping onto the description I outlined above, in terms of Eq. (2.1) and Eq. (2.2). The Schrödinger equation in Eq. (2.3) is just a differential equation with \hbar appearing as a parameter. You might have thought that one would expand ψ in a Taylor series in \hbar like,

What you cannot do

$$\psi = \psi_0 + \hbar\psi_1 + \hbar^2\psi_2 + \dots, \quad (2.5)$$

plug it into the equation and try to solve it order by order in \hbar . The $\psi_0, \psi_1 \dots$ will all have weird dimensions since \hbar is not dimensionless; this, however, is not a serious issue. The key point is that, in such an expansion, we are assuming ψ to be analytic in \hbar . This Taylor series expansion, however, does not work, as you can easily verify. In fact, we would have been in a bit of trouble if it had worked since we would then have to interpret ψ_0 as some kind of “classical” wavefunction. The way one obtains the

classical limit is quite different. We will get it from the ansatz in Eq. (2.1) which has \hbar occurring *non-analytically* in the phase.

What you can do

Let us now carry out the procedure described earlier. Using the expression for ψ from Eq. (2.1) in Eq. (2.3), and equating the real and imaginary parts, we get the two equations

$$(R^2 S')' = -m \frac{\partial R^2}{\partial t} \quad (2.6)$$

and

$$\frac{S'^2}{2m} + V(q) + \frac{\partial S}{\partial t} = \frac{\hbar^2}{2m} \frac{R''}{R}, \quad (2.7)$$

*Same math: but
some physics is lost*

where the prime denotes the derivative with respect to q . The Schrödinger equation is completely equivalent to the two real equations in Eq. (2.6) and Eq. (2.7). Anything you can do with a complex wavefunction ψ can also be done with two real functions R and S . But, of course, the Schrödinger equation is linear in ψ while equations Eq. (2.6) and Eq. (2.7) are non-linear, thereby hiding the principle of superposition of quantum states — which is a cornerstone of the quantum description.

Equation (2.7) suggests an alternate scheme for doing the Taylor series expansion in \hbar . We can now try to interpret the left hand side of Eq. (2.7) as the lowest order contribution to the phase of the wavefunction in Eq. (2.1). In such a case, we can attempt a Taylor series expansion in the form

$$S(t, q) = S_0(t, q) + \hbar^2 S_1(t, q) + \dots \quad (2.8)$$

This means the leading behaviour of the wavefunction is given by $\exp(iS_0/\hbar)$ which is *non-analytic* in \hbar . This is a different kettle of fish when it comes to a series expansion in terms of a parameter in a differential equation. Also, note that Eq. (2.7) depends only on \hbar^2 and not on \hbar ; so the second term in the Taylor series starts with \hbar^2 , and not with \hbar .

*Quantum to
Classical = Wave
optics to Ray optics*

Why does this approach work while the expansion in Eq. (2.5) does not lead to sensible results? The reason essentially has to do with the fact that — in proceeding from quantum physics to classical physics — we are doing something analogous to obtaining ray optics from electromagnetic waves. One knows that this can come about only when the phase of the wave is non-analytic in the expansion parameter — which is essentially the wavelength in the case of light propagation. So you need to bring in some *extra physical insight* to obtain the correct limit.

While ψ is non-analytic in \hbar , we have now translated the problem into R and S which are (assumed to be) analytic in \hbar so that the standard procedure works. To the leading order, we will ignore the right hand side of Eq. (2.7) and obtain the equation

$$\frac{S_0'^2}{2m} + V(q) + \frac{\partial S_0}{\partial t} = 0. \quad (2.9)$$

(This result might seem obvious but there is a subtlety lurking here which we will comment on later.) This partial differential equation determines the phase of the wavefunction to the lowest order of accuracy in \hbar . Solving it is pretty easy; you try an ansatz $S_0(t, q) = -(t - t_0)E + F(q)$ where E and t_0 are two constants. An elementary integration gives the solution as

$$S_E(t, q) = -(t - t_0)E + \int dq \sqrt{2m(E - V(q))} , \quad (2.10)$$

which depends on E as a parameter, indicated explicitly by a subscript in S_E . (I have dropped the subscript “0” for simplicity.) Strictly speaking, the square root in Eq. (2.10) comes with a \pm factor in front; we have chosen one of the branches using an initial condition on the direction of the velocity. Correspondingly, the wavefunction is given by

$$\psi_E(t, q) \simeq R \exp \frac{1}{\hbar} \left[-iE(t - t_0) + i \int dq \sqrt{2m(E - V(q))} \right] , \quad (2.11)$$

which again depends on E as a parameter. So far, we have merely written down the Schrödinger equation, Eq. (2.3), and solved it in a particular approximation. Where is classical physics and where are the trajectories?

To obtain the classical trajectory out of this quantum wavefunction, we use the idea of constructive interference of waves. Since E is just a parameter and the Schrödinger equation in Eq. (2.3) is linear in ψ , we can superpose solutions with different values of E to construct a wave packet. When we add ψ_E with different values of E , the condition for constructive interference corresponds to the phase of the wavefunction remaining stationary when E changes by a small amount ΔE . That is, we impose the condition

$$S_E(t, q) = S_{E+\Delta E}(t, q) . \quad (2.12)$$

This is equivalent to the condition $(\partial S_E / \partial E) = 0$. For S_E in Eq. (2.10), this leads to

$$t - t_0 = \int dq \left(\frac{m}{2(E - V)} \right)^{1/2} , \quad (2.13)$$

which gives you the sought-after trajectory $q(t)$ as a function of the parameter E . (You need to fix the two parameters E and t_0 by the boundary conditions of the problem.) The phase of the wavefunction singles out this trajectory in the $t - q$ plane by the condition of constructive interference in Eq. (2.12). The explicit emergence of the classical trajectory is shown graphically in Fig. 2.1 for the simple potential $V(x) = mgx$.

It is elementary to show from Eq. (2.13) that the trajectory satisfies the equations Eq. (2.4) with $H(p, q) = E$. The second of these equations (viz. Newton’s second law) is *not* the most efficient way to solve for the trajectory of the particle. Almost always, solving Eq. (2.9) and demanding

Warning! See later!

Condition for constructive interference

The magical emergence of classical trajectory!

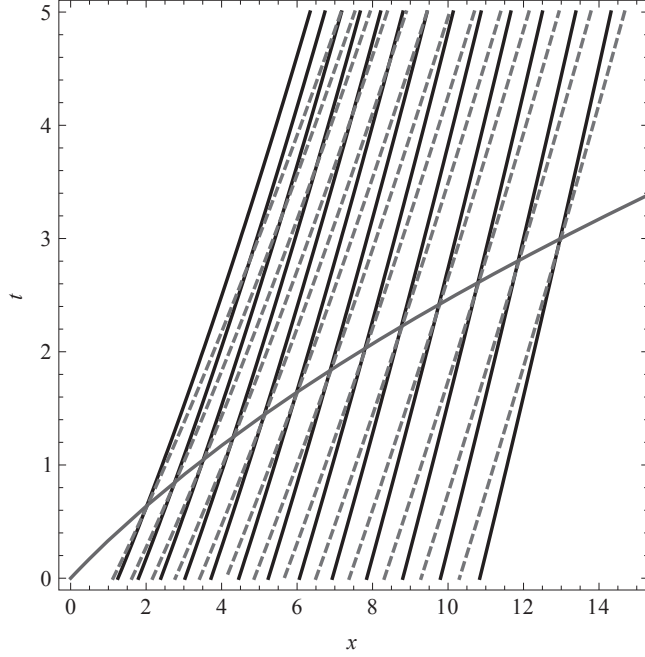


Fig. 2.1: The emergence of a classical trajectory from the constructive interference of quantum phases. As an illustration, we consider curves of constant phase $S_E(t, x)$ in the $t - x$ plane for the energies E and $E + \Delta E$. The function $S_E(t, x)$ is evaluated using Eq. (2.10) for the potential $V = mgx$. The set of unbroken curves are given by $S_E(t, x) = \text{constant}$, while the dashed curves are for $S_{E+\Delta E}(t, x) = \text{constant}$. The condition of constructive interference requires S to remain unchanged when $E \rightarrow E + \Delta E$. This condition $S_E(t, x) = S_{E+\Delta E}(t, x)$ determines a set of points on the $t - x$ plane shown in the figure which passes through the intersection points of the two families of curves. This is the classical trajectory given by $x = (1/2)gt^2$ with suitable initial conditions.

$(\partial S / \partial E) = 0$ is a faster route to the trajectory. Quantum physics gives the most efficient route to the classical trajectory!

Let us pause and savour what we have achieved. We started with the Schrödinger equation for a particle in a potential V and determined the phase of the wavefunction to the lowest order of accuracy in \hbar . This phase satisfied a partial differential equation, Eq. (2.9). The solution to this partial differential equation introduced the parameter E into the problem so that the phase of the wavefunction depended on this parameter E . We then looked for the region in the $t - q$ plane in which constructive interference of the waves, with different values of E occurs. This is equivalent to demanding $(\partial S / \partial E) = 0$ and it singled out the trajectory followed by the particle in the $t - q$ plane.

If we want to forget about quantum mechanics and only want to know the classical trajectory of a particle in a potential V , then we can express the whole procedure in an algorithmic fashion:

1. Define a Hamiltonian $H(p, q)$. In our case, it was $H(p, q) = (p^2/2m) + V(q)$ but it could have been more general.
2. Write down the partial differential equation for a function $S(t, q)$ given by

$$\frac{\partial S}{\partial t} + H\left[\frac{\partial S}{\partial q}, q\right] = 0 \quad (2.14)$$

which arises as the lowest order approximation to the equation satisfied by the wavefunction. Solve this partial differential equation, which will introduce the constants E and t_0 leading to the solution $S(t, q; E, t_0)$.

This function is called the *action* purely because of historical reasons.

3. Impose the condition $(\partial S/\partial E) = 0$. This will give you the classical trajectory taken by the particle in terms of the two arbitrary constants E and t_0 . Fix the constants using the boundary conditions of the problem.

You might recognize Eq. (2.14) as the Hamilton-Jacobi equation from a classical mechanics course. It most probably was introduced after a lot of talk about the so-called canonical transformations, generating functions and what not. The condition $(\partial S/\partial E) = 0$ would have come as a condition on new coordinates and momenta in a canonical transformation. Forget it all! Particles do not follow trajectories. They are described by wavefunctions but under appropriate circumstances the constructive interference of the phases of the wavefunction will single out a path which we call a classical trajectory. The Hamilton-Jacobi equation is just the lowest-order Schrödinger equation if we use the ansatz in Eq. (2.1). The mysterious procedure in Hamilton-Jacobi theory — of differentiating the solution to Hamilton-Jacobi equation and equating it to a constant — is just the condition for constructive interference of the phases of waves differing slightly in the parameter E . The procedure based on Hamilton-Jacobi theory works in classical mechanics because it is supported by the Schrödinger equation.

If you haven't seen the Hamilton-Jacobi equation before, nothing is lost!

Box 2.1: The Hamilton-Jacobi equation is a dispersion relation!

The Hamilton-Jacobi equation is essentially a dispersion relation for a complex wave. This is easy to see in the context of non-relativistic quantum mechanics. If a quantum amplitude is expressed in the form $\psi = R \exp(iS/\hbar)$, then the Hamilton-Jacobi equation relates $p = \partial S/\partial q$ to $E = -\partial S/\partial t$ by the condition $p^2(q) = 2m(E - V)$. This is a relation between the wave vector $k = p/\hbar$ and the frequency $\omega = E/\hbar$ of the “matter wave” associated with the particle.

In fact, this idea generalizes to the relativistic case as well. In this case, the Schrödinger equation will be replaced by a more complicated equation, say, the Klein-Gordon equation, which might also include interaction terms with electromagnetic or gravitational fields.

*Relativity before
relativity: $H = -\partial_t S$
and $\mathbf{p} = \nabla S$ is
just $p_a = \partial_a S$*

Though the probabilistic interpretation will no longer hold for the solutions in general, it can be made to work in the appropriate limit and the classical trajectory can still be obtained by the same prescription as in non-relativistic quantum mechanics. We again express the solutions to the relevant wave equation in the form $\Psi = R \exp(iS/\hbar)$ and define the *four-momentum* of the particle as $p_a = \partial_a S$, which nicely incorporates the results $(\partial S/\partial t) = -E$, $\nabla S = \mathbf{p}$ at one go. The Hamilton-Jacobi equation can be now obtained from the known relation between the energy and momentum.

For example, a free relativistic particle has $\eta^{ij} p_i p_j = -m^2 c^2$ which is just a fancy way of writing the relation between energy and momentum: $E^2 = |\mathbf{p}|^2 c^2 + m^2 c^4$. The Hamilton-Jacobi equation is obtained by replacing p_j by $\partial_j S$ to give: $\eta^{jk} \partial_j S \partial_k S = -m^2 c^2$.

*Same story for a
particle in an
electromagnetic
field ...*

A more non-trivial case is a charged particle in an electromagnetic field described by a vector potential A_i . In this case, the four-momentum changes as: $p_j \rightarrow (p_j - qA_j)$. The corresponding Hamilton-Jacobi equation is:

$$\eta^{jk} (\partial_j S - qA_j) (\partial_k S - qA_k) = -m^2 c^2 . \quad (2.15)$$

If you solve this equation in a given electromagnetic potential A_k and impose the condition for constructive interference, you will get the trajectory of the charged particle in this field. (We will see an example in Chapter 3.)

The situation with the gravitational field is even simpler. Gravity is described by changing the special relativistic line interval $ds^2 = \eta_{ij} dx^i dx^j$ to the form $ds^2 = g_{ij} dx^i dx^j$, where g_{ij} is the metric tensor which describes the curved spacetime and gravity. (You will learn why, in Chapter 11.) The dispersion relation for momentum now changes from $\eta^{ab} p_a p_b = -m^2 c^2$ to $g^{ab} p_a p_b = -m^2 c^2$. Substituting $p_a = \partial_a S$ then gives you the Hamilton-Jacobi equation in the presence of gravity

*.... and for a
particle in a
gravitational field.*

$$g^{ab} \partial_a S \partial_b S = -m^2 c^2 . \quad (2.16)$$

The rest of the algorithm to get the trajectory is the same as before. The equations, (2.15), (2.16) etc. describe the dispersion relations for waves associated with material particles interacting with electromagnetic or gravitational fields in the $\hbar \rightarrow 0$ limit.

As I explained at the beginning of this chapter, the ideas developed here are extremely general and — in fact — *we do not know of any physical system which is not encompassed by these principles.*

This looks good, but haven't we overstepped our limits? Surely there must exist quantum states described by some ψ which do not lead to classical trajectories? What happened to them? Sure there are; to see where they fit in, let us study couple of examples.

*Quantum states
without classical
limit*

To begin with, note that, though we developed the above approach from a desire to obtain the classical limit, mathematically speaking, we are just studying an approximation to the differential equations governing the system — usually known as Wentzel-Kramers-Brillouin (WKB) approximation. This fact is strikingly evident in the context of quantum mechanical tunneling which, of course, has no classical analogue. Nevertheless, we can get a reasonable approximation to the wavefunction in a classically forbidden form by taking $E < V(q)$ in Eq. (2.9). In this range, say, $a < q < b$ where $E < V(q)$, we see that S_0 picks up the imaginary part given by

*Example 1:
Tunneling*

$$S_0 = \sqrt{2m} \int_a^b \sqrt{E - V(q)} dq = i\sqrt{2m} \int_a^b \sqrt{V(q) - E} dq . \quad (2.17)$$

The wavefunction now becomes exponentially decreasing (or increasing) in this classically forbidden range. Without the oscillatory behaviour, so there is no constructive interference of waves and no classical trajectories!

The second context is related to the subtlety which I mentioned earlier in ignoring the right hand side of Eq. (2.7). For this approximation to be valid, we must have

*Example 2:
Ground states*

$$\lim_{\hbar \rightarrow 0} \frac{\hbar^2}{2m} \frac{R''}{R} = 0 . \quad (2.18)$$

It is easy to construct states for which this condition is violated! As a simple example consider the ground state of a system in a bounded potential which will be described by a real wavefunction. In this case, $\psi = R$ and $S = 0$. From Eq. (2.7) we now see that

$$\frac{\hbar^2}{2m} \frac{R''}{R} = V(q) - E . \quad (2.19)$$

The limit in Eq. (2.18) cannot now hold, in general. Clearly our analysis fails for the ground state of a quantum system when we try the ansatz in Eq. (2.1). To see this explicitly, consider the ground state of a harmonic oscillator:

$$\psi(q) = N \exp \left[-\frac{m\omega}{2\hbar} q^2 \right] . \quad (2.20)$$

This wavefunction is an exact solution to the Schrödinger equation, and its amplitude and phase (which is zero) must thus satisfy Eq. (2.6) and Eq. (2.7). A straightforward computation now shows — not surprisingly — that

$$\frac{\hbar^2}{2m} \frac{R''}{R} = \frac{1}{2} m \omega^2 q^2 - \frac{1}{2} \hbar \omega . \quad (2.21)$$

When we take the limit $\hbar \rightarrow 0$, the second term on the right hand side vanishes but not the first term! This means that there are quantum states for which we cannot naively take the right hand side of Eq. (2.7) to be zero and determine the classical limit. (Interestingly enough, this is also true for the time-dependent, coherent, states of the oscillator. You may want to amuse yourself by analyzing this situation in greater detail.) The $\hbar \rightarrow 0$ limit of the Schrödinger equation is far from trivial.

Box 2.2: The Wigner function

Can we interpret the wavefunction *itself* in the classical limit rather than obtain a trajectory by constructive interference and use it to describe classical limit? This is tricky and the best we can do is to use the concept of the Wigner function $F(q, p, t)$, corresponding to a wavefunction $\psi(q, t)$, defined by the relation

$$F(q, p, t) = \int_{-\infty}^{\infty} du \psi^* \left(q - \frac{1}{2} \hbar u, t \right) e^{-ipu} \psi \left(q + \frac{1}{2} \hbar u, t \right). \quad (2.22)$$

The idea is to see whether one can think of F as a probability distribution function in the phase space (with position (q) and momentum (p) as coordinates) since F simultaneously encodes both coordinate space and momentum space information in a state represented by ψ . (Some of the pedagogical details regarding Wigner functions can be found, e.g., in Ref. [3].) If you integrate F over the momentum variable p , you get

$$\int_{-\infty}^{\infty} dp F(q, p, t) = |\psi(q, t)|^2, \quad (2.23)$$

while if you integrate F over q you get

$$\int_{-\infty}^{\infty} dq F(q, p, t) = |\phi(p, t)|^2, \quad (2.24)$$

where $\phi(p, t)$ is the Fourier transform of $\psi(q, t)$. From the standard rules of quantum mechanics, we know that $\phi(p, t)$ gives the probability distribution in the momentum space. Therefore, when marginalized over either coordinate, F satisfies the nice properties that we expect of a probability distribution. Further, direct differentiation of Eq. (2.22) and some clever manipulation will allow you to obtain the following equation satisfied by F :

$$\frac{\partial F}{\partial t} + \frac{p}{m} \frac{\partial F}{\partial q} - \frac{dV}{dq} \frac{\partial F}{\partial p} = \frac{\hbar^2}{24} \frac{d^3 V}{dq^3} \frac{\partial^3 F}{\partial p^3} + \dots, \quad (2.25)$$

So far, so good!

where \dots denotes terms which are of higher order in \hbar . So if the potential is at most quadratic in the coordinates, or for *arbitrary* potentials up to linear order in \hbar , the right hand side of Eq. (2.25) vanishes and we get exactly the continuity equation in phase space with the semi-classical identifications $\dot{q} = p/m$ and $\dot{p} = -V'$.

The only trouble — and a serious one — is that F is not positive-definite and since we do not know how to interpret negative probabilities, we cannot use F as a probability distribution in the phase space. For example, the Wigner function corresponding to the first excited state of a harmonic oscillator (in suitable units) is

$$F(q, p) = 4(q^2 + p^2 - (1/2))e^{-(p^2 + q^2)}, \quad (2.26)$$

which can go negative.

This does not, however, prevent us from using the Wigner function in suitable limits as an *approximation* to the classical probability distribution. In particular, the Wigner function corresponding to the semi-classical wavefunction is quite easy to interpret. In this case, we get

$$F(q, p) \propto \frac{1}{|S'_0(q)|} \delta_D \left(p - \frac{\partial S_0}{\partial q} \right) + \mathcal{O}(\hbar^2). \quad (2.27)$$

The Dirac delta function δ_D tells you that when the particle is at q , its momentum is sharply peaked at $(\partial S_0/\partial q)$ which is exactly what we would have expected if the particle was moving along a classical trajectory. Further, the probability to find the particle around q is proportional to $(1/S'(q))$ which is in turn proportional to the time $dt = dq/v(q)$ (where $v(q)$ is the speed of the particle) which the particle spends in the interval $(q, q + dq)$. Note that now the probability distribution is *not* peaked around any single trajectory $q(t)$; however, once you pick a q , it gives you a unique p . This correlation between momentum and position is the key feature of the classical limit.

If you take a classically forbidden region (in which the wavefunction is exponentially damped, rather than oscillatory), which is a “purely quantum mechanical” situation, you will find that the Wigner function factorizes into a product of two functions: $F(q, p) = F_1(q)F_2(p)$. The momentum and position are totally uncorrelated in such a state which clearly is the other extreme of the semi-classical state in which the momentum is completely correlated with the position. The same decoupling of momentum and position dependence occurs for many other states. One simple example is the ground state of the harmonic oscillator for which you will find that the Wigner function factorizes into two products, both Gaussian, in position and momentum. So the ground state of the harmonic oscillator is as non-classical as a state could get in this interpretation.

The fly in the ointment

Works for the semiclassical states

We can express the action S in a different form which turns out to be extremely valuable. To do this, we note that Eq. (2.10) can be expressed as an integral over time as

$$S(t, q) = - \int_{t_0}^t dt E + \int_{t_0}^t \dot{q} dt \sqrt{2m(E - V)}. \quad (2.28)$$

Along the classical trajectory determined by Eq. (2.4) with $E = (1/2)m\dot{q}^2 + V$, this becomes:

$$\begin{aligned} S(t, q) &= \int_{t_0}^t dt \left[- \left(\frac{1}{2}m\dot{q}^2 + V \right) + m\dot{q}^2 \right]_C \\ &= \int_{t_0}^t dt \left[\frac{1}{2}m\dot{q}^2 - V \right]_C. \end{aligned} \quad (2.29)$$

The subscript C is a reminder that the integral is evaluated along the classical trajectory connecting, say, some q_0 at $t = t_0$ with q at time t . The resulting $S(t, q)$ is treated as a function of (t, q) . We will now treat this equation as defining S for *all* paths $q(t)$ which start at some fixed q_0 at time t_0 but are otherwise completely arbitrary; that is, the functions $q(t)$ need not necessarily satisfy Eq. (2.4). So we now consider the quantity $S[t, q; q(t)]$ defined by

The most useful way to write S

$$S(t, q; q(t)) = \int_{t_0, q_0}^{t, q} dt L(\dot{q}, q); \quad L \equiv \frac{1}{2}m\dot{q}^2 - V(q), \quad (2.30)$$

in which S depends on the upper end point of integration as well as on the path $q(t)$. Let us now vary the end point from (t, q) to $(t + \delta t, q + \delta q)$. Originally, suppose we had a path $q(t)$ which connected (t_0, q_0) and (t, q) . After varying the end points, we will have a *different* path $q(t) + \delta q(t)$ which connect (t_0, q_0) and $(t + \delta t, q + \delta q)$. We can compute the value of S for both these trajectories and ask what is the change δS due to our change in the end point.

There are two ways of computing this change δS . The straightforward way of computing δS is to treat it as a function of q and t at the end point and evaluate the change in terms of partial derivatives using Eq. (2.14):

$$\delta S = \frac{\partial S}{\partial q} \delta q + \frac{\partial S}{\partial t} \delta t = p \delta q - H \delta t. \quad (2.31)$$

There is, however, another way of computing it by explicitly varying the path $q(t)$, as well as the end points, in the expression for S in Eq. (2.30). This will lead to the result

$$\delta S = L(q, \dot{q}) \delta t + \int_{q_0, t_0}^{q + \Delta q, t} \delta L dt. \quad (2.32)$$

The first term arises because we changed t to $t + \delta t$ at the end point. The second term has two contributions: (a) When the path is changed, L changes by the amount:

$$\delta L = \frac{\partial L}{\partial q} \delta q + \frac{\partial L}{\partial \dot{q}} \delta \dot{q} = \left[\frac{\partial L}{\partial q} - \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}} \right) \right] \delta q + \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}} \delta q \right) \quad (2.33)$$

and (b) the end point of the path changes from q to $q + \Delta q$. This allows us to write Eq. (2.32), on using $p = (\partial L / \partial \dot{q})$, as:

$$\delta S = \int dt \left(\frac{\partial L}{\partial q} - \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}} \right) \right) \delta q + p \Delta q + L \delta t. \quad (2.34)$$

The Δq , in turn, is made of two pieces. First, there is an ‘intrinsic’ change due to δq at the end points. Second, when one makes only a δt change in the end point, one induces a change $(-\dot{q} \delta t)$ in q (see Fig. 2.2). Hence, the total change in q at constant t is given by $\Delta q = \delta q - \dot{q} \delta t$. Using $\Delta q = \delta q - \dot{q} \delta t$ in Eq. (2.34), we get:

This is non-trivial; see Fig. 2.2!

$$\delta S = \int dt \left(\frac{\partial L}{\partial q} - \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}} \right) \right) \delta q + p \delta q + (L - p \dot{q}) \delta t. \quad (2.35)$$

Equating the expressions for δS in Eq. (2.35) and Eq. (2.31) and recalling that $H = p \dot{q} - L$, we find that the classical trajectory must satisfy the condition

$$\frac{\partial L}{\partial q} - \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}} \right) = 0. \quad (2.36)$$

In other words, one can also determine the classical trajectory by starting from the definition of action in Eq. (2.30) and demanding that $\delta S = 0$ for variations of the trajectories with $\delta q = 0$ at the end points.

This gives us the standard Lagrangian-based action principle in classical mechanics. But note that such a variational principle means nothing within the context of classical theory! Classically, a particle is supposed to follow a *specific trajectory* and — at best — the value of S for this classical trajectory could have some meaning. Defining a quantity S for an arbitrary trajectory has no physical meaning within a classical theory. The situation is quite different in quantum mechanics in which we have no unique trajectory at all. Rather, all possible trajectories co-exist and a classical trajectory is selected by the constructive interference condition. We used the condition previously to give meaning to the Hamilton-Jacobi equation. But one can play the same game with the action principle itself, which provides a powerful technique in quantum theory. We shall say more about this in Chapter 17.

Quantum mechanics gives meaning to classical action principle

You would have noticed that we have actually proceeded in a direction *opposite* to the standard textbook description! Normally, you would have started in your physics course with an action principle based on a

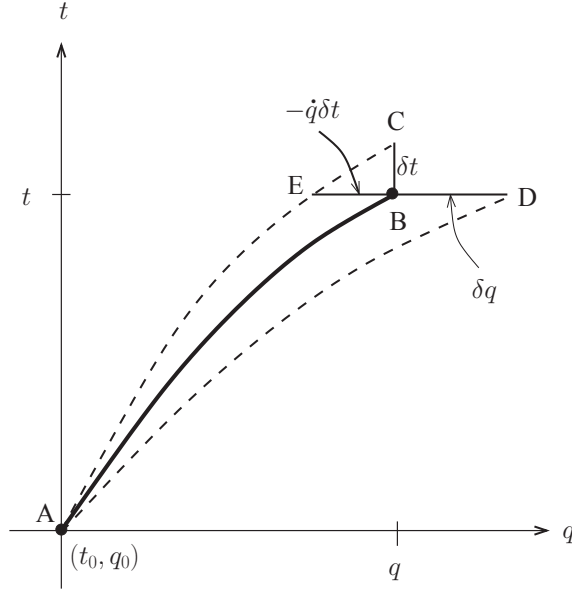


Fig. 2.2: The variation of the action during the change of end points. The original path connects the event $A(t_0, q_0)$ with the event $B(t, q)$. Varying t to $t + \delta t$ with fixed q , shifts B to C (and the trajectory shifts from AB to AC). Varying q to $q + \delta q$ with fixed t , shifts B to D (and the trajectory shifts from AB to AD). Under simultaneous variation, the change in q consists of two parts: (i) The part BD corresponding to the intrinsic change δq , and (ii) the part BE induced by the change δt in t , given by $(-\dot{q}\delta t)$. Therefore, the net change in q at constant t is given by $\Delta q = \delta q - \dot{q}\delta t$.

Lagrangian $L(\dot{q}, q)$, defined a Hamiltonian as $H = p\dot{q} - L$, written down a Hamilton-Jacobi equation as in Eq. (2.14), etc. The Schrödinger equation, on the other hand, leads naturally to the Hamilton-Jacobi equation and the functional form of $S(t, q)$ obtained by integrating the Lagrangian *for the classical trajectory*. So when you go from quantum mechanics to classical mechanics, the usual procedure is indeed reversed!

As an aside, let me also comment on another issue related to the form of the action functional that we have obtained in Eq. (2.30). The action which we found directly from quantum theory, as the phase of the wavefunction, has the form of Eq. (2.10). If we leave out the time dependence of the phase, we are left with an action which is an integral of $p(q, E)dq$ where $p(q, E)$ is the momentum of the particle with energy E when it is at the location q . This is structurally quite different from the action in Eq. (2.30) and it is worth analyzing it a bit.

To see what is involved, let us generalize from one dimension to, say, $D = 3$ and consider a situation in which we are only interested in various trajectories connecting two points in space x_1^α and x_2^α , irrespective of their parameterization. We can then describe the curves with some param-

eter λ by giving the D functions $x^\alpha(\lambda)$. This parameter λ is irrelevant and we could have described the same curves by some other set of functions obtained by changing the parameterization from $\lambda \rightarrow f(\lambda)$. We now want to construct an extremum principle from which we can determine the classical trajectory as a geometrical curve without any reference to the parametrization. The relevant action functional for this extremum problem (called the Jacobi-Mapertuis action) is given by

$$A_J = \int_{x_1}^{x_2} p_\alpha dx^\alpha = \int_{\lambda_1}^{\lambda_2} d\lambda \left(\frac{\partial L}{\partial \dot{x}^\alpha} \right) \dot{x}^\alpha, \quad (2.37)$$

where the overdot now represents derivative with respect to the parameter λ used to parametrize the curves $x^\alpha(\lambda)$. It can be easily verified that the modified action principle based on A_J in Eq. (2.37) leads to the actual *paths in space* as a solution to the variational principle $\delta A_J = 0$ when we vary all trajectories connecting the end points x_1^α and x_2^α . This trajectory will have energy E , but will contain no information about the time coordinate. In fact, the first form of the action in Eq. (2.37) makes clear that there is no time dependence in the action. The curves remain the same and the invariance of the action under the reparameterization expresses this fact.

Path finder

It is possible to rewrite the expression for A_J in a nicer form. Using the fact that L is a homogeneous quadratic function of velocities, we have the result

$$\dot{x}^\alpha \left(\frac{\partial L}{\partial \dot{x}^\alpha} \right) = 2T = m \left(\frac{d\ell}{d\lambda} \right)^2 = 2[E - V(x^\alpha)], \quad (2.38)$$

where T is the kinetic energy and ℓ is the arc length of the path. Substituting into Eq. (2.37) we get

$$A_J = \int_{x_1}^{x_2} m \left(\frac{d\ell}{d\lambda} \right) d\ell = \int_{x_1}^{x_2} \sqrt{2m(E - V(x^\alpha))} d\ell, \quad (2.39)$$

which is again manifestly reparameterization invariant with no reference to time. In some sense, this is a more natural form of the action which arises directly from quantum theory rather than the action in Eq. (2.30). After all, we are interested in the classical trajectory, not how that trajectory is parameterized. We will say more about this action in Chapter 17.

All this is fine as long as you are told that the Hamiltonian is $H = p^2/2m + V$ (so that you can write the Schrödinger equation) or that the Lagrangian is $L = (1/2)m\dot{q}^2 - V$ (so that you can develop quantum mechanics from an action principle; see Chapter 17). *But how do we know this?* Why is the Lagrangian given by such a strange combination? The reason is pretty non-trivial and illustrates the point that *exact theories make more sense than approximate ones*.

But why is the Lagrangian $K - V$?

Let us first consider the non-relativistic free particle for which the action is an integral over the kinetic energy. It is not very clear why min-

Action makes better sense in special relativity than in non-relativistic mechanics

imizing this quantity should have any physical significance. But let us next consider the *special relativistic* free particle following a worldline in spacetime along some arbitrary curve with speed $v(t)$. We attach a clock to the particle and ask how much time ($\Delta\tau$) will elapse in this moving clock, when a stationary clock in the lab frame \mathcal{S} shows a lapse of Δt . At any instant t , the particle is momentarily at rest in a comoving Lorentz frame (\mathcal{S}') boosted with respect to \mathcal{S} by some velocity $\mathbf{v}(t)$. Since the interval $ds^2 = -c^2 dt^2 + d\mathbf{x}^2$ has the same value in all Lorentz frames, we can evaluate it in \mathcal{S} and (\mathcal{S}') and equate the results. In the comoving frame of the clock (\mathcal{S}'), we have $ds^2 = -c^2 d\tau^2$ since $d\mathbf{x}^2 = 0$, while in \mathcal{S} , we have $ds^2 = -c^2 dt^2 + d\mathbf{x}^2 = -c^2 dt^2 [1 - v^2(t)/c^2]$. So we get:

$$\tau = \int dt \left(1 - \frac{v^2(t)}{c^2} \right)^{1/2}, \quad (2.40)$$

which — called the proper time — is clearly an invariant quantity. Note that this expression is valid for clocks in an arbitrary state of motion, including accelerated motion. (I stress this because students sometimes think that this result is valid only for inertial motion of the clock.)

It makes some physical sense to claim that ‘particles follow a trajectory of least time’ and take the action to be proportional to τ . If we take the proportionality constant as $-mc^2$, we can ensure a suitable limit when $(v/c) \ll 1$:

$$\begin{aligned} S &= -mc^2 \tau \\ &= -mc^2 \int dt \left(1 - \frac{v^2(t)}{c^2} \right)^{1/2} \rightarrow \int dt \left[\frac{1}{2}mv^2 - mc^2 \right]. \end{aligned} \quad (2.41)$$

So you see, the action for a *non-relativistic* particle, being an integral over kinetic energy, acquires the nice interpretation of extremizing the proper time, *in special relativity*, if we ignore the constant $-mc^2$.

Warning! See Chapter 15.

This is fine for a free particle but what about a particle in an electromagnetic or a gravitational field? We now have to make sure that any external field we introduce respects special relativity. This limits the kind of expressions we can integrate over to get S . We can only use

$$S = c_1 \int ds + c_2 \int A_j dx^j + c_3 \int \sqrt{-g_{ij} dx^i dx^j}, \quad (2.42)$$

up to quadratic order, where the c_i -s are constants, A_j is a four-vector and g_{ij} is a second rank tensor. Since $ds = \sqrt{-\eta_{ij} dx^i dx^j}$, you can get the first term as a special case of the last by taking $g_{ij} = \eta_{ij}$. So, up to quadratic order, we can *only* use

$$S = c_2 \int A_j dx^j + c_3 \int \sqrt{-g_{ij} dx^i dx^j}, \quad (2.43)$$

The raison d'être for just two classical fields

with just two external fields: A_j gives you electromagnetism and g_{ij} gives you gravity!

With this structure, it is easy to show that in *non-relativistic electrostatics* the Lagrangian will turn out to be kinetic energy *minus* electrostatic potential energy; this is *not* true in general — even in electromagnetism it is not true when we go beyond electrostatics. The reason we use kinetic energy *minus* potential energy for gravity is a lot more beautiful. Believe it or not, this is because gravity affects the flow of time!! We will learn this in Chapter 11.

Orbits of Planets are Circles!

3

Yes, it is true. And no, it is not the cheap trick of tilting the paper at an angle to see an ellipse as a circle. The real trick is a bit more sophisticated. It turns out that the trajectory of a particle, moving under the attractive inverse square law force, is a circle (or part of a circle) in the *velocity space* (The high-tech name for the path in velocity space is *hodograph*).

The proof is quite trivial and, in fact, the entire Kepler problem is quite trivial but for the textbooks making it complicated. If you think straight you can solve it in couple of steps, as I will now show.

You would have solved Kepler problem in two steps if only you haven't taken courses in classical mechanics!

To set the stage, we start with the result that, for particles moving under *any* central force $f(r)\mathbf{e}_r$, the angular momentum $\mathbf{J} = \mathbf{r} \times \mathbf{p}$ is conserved. Here \mathbf{r} is the position vector, \mathbf{p} is the linear momentum and \mathbf{e}_r is the unit vector in the direction of \mathbf{r} . This implies that the motion is confined to the plane perpendicular to \mathbf{J} which we take to be $\theta = \pi/2$ plane. (The constancy of $J = mr^2\dot{\phi}$ also gives Kepler's second law, since $(r^2\dot{\phi}/2) = J/2m \equiv h/2$ is the area swept by the radius vector in unit time.) Let us introduce in this plane the polar coordinates (r, ϕ) as well as the Cartesian coordinates (x, y) . Let the unit vector in the ϕ direction be \mathbf{e}_ϕ with Cartesian components $(-\sin \phi, \cos \phi)$ which satisfies the relation $d\mathbf{e}_\phi/d\phi = -\mathbf{e}_r$. So we have (this is the first of the two-step derivation!):

Step 1: Solve for \mathbf{v}

$$\frac{d\mathbf{v}}{d\phi} = \frac{\dot{\mathbf{v}}}{\dot{\phi}} = -\frac{GM}{h}\mathbf{e}_r = \frac{GM}{h}\frac{d\mathbf{e}_\phi}{d\phi}, \quad (3.1)$$

where we have used $r^2\dot{\phi} \equiv h$ and $\dot{\mathbf{v}} = -(GM/r^2)\mathbf{e}_r$ to arrive at the second equality. It follows that $\mathbf{v} - (GM/h)\mathbf{e}_\phi$ is a constant vector which we will denote by \mathbf{w} . Therefore

$$\mathbf{v} = \mathbf{w} + \frac{GM}{h}\mathbf{e}_\phi. \quad (3.2)$$

Step 2: Get the trajectory!

Taking a dot product of this equation with \mathbf{e}_ϕ we obtain (this is the second and final step of the derivation!)

$$\mathbf{v} \cdot \mathbf{e}_\phi = v_\phi = \frac{h}{r} = w \cos \phi + \frac{GM}{h} \equiv \frac{GM}{h} (1 + e \cos \phi), \quad (3.3)$$

where we have used $v_\phi = r\dot{\phi} = h/r$ and defined constant e by the relation $w \equiv (GM/h)e$. This is clearly a conic section with eccentricity e and latus rectum (GM/h^2) . We have also oriented the axes such that \mathbf{w} is along the y-axis so that the angle ϕ between \mathbf{w} and \mathbf{e}_ϕ is the same as the angle between \mathbf{r} and the x-axis. So you see, it is really easy.

Most of the remaining part of the chapter is devoted to appreciating different aspects of this result and we will do it slowly, savoring every moment.

Let us start with the result in Eq. (3.2), which tells you that $(\mathbf{v} - \mathbf{w})^2 = (GM/h)^2$; that is, the tip of the vector \mathbf{v} moves on a circle of radius GM/h centered at \mathbf{w} ! To see this more explicitly, let us choose (say) $v_x(\phi = 0) = 0$; $v_y(\phi = 0) = u$ and obtain from Eq. (3.2) the result:

$$\mathbf{w} = \frac{GM}{h} e \hat{\mathbf{y}} = \left[u - \left(\frac{GM}{h} \right) \right] \hat{\mathbf{y}} \quad (3.4)$$

so that $u = (GM/h)(1 + e)$. (Here $\hat{\mathbf{y}}$ is a unit vector in the y-direction.) Then \mathbf{v} satisfies the condition (the “hodograph”):

$$v_x^2 + \left(v_y - \frac{GM}{h} e \right)^2 = \left(\frac{GM}{h} \right)^2, \quad (3.5)$$

which is a circle with center at $(0, eGM/h)$ and radius GM/h . So you see, planets do move in circles, as advertised!

It is clear that the structure of the hodograph depends vitally on the ratio between u and GM/h ; that is on e . The geometrical meaning of e is clear from Fig. 3.1.

All standard results recovered

- If $e = 0$, i.e., if we had chosen the initial conditions such that $u = GM/h$, then the center of the hodograph is at the origin of the velocity space and the magnitude of the velocity remains constant. Writing $h = ur$, we get $u^2 = GM/r$ leading to a circular orbit in the real space as well.
- When $0 < e < 1$, the origin of the velocity space is inside the hodograph. As the particle moves, the magnitude of the velocity changes between a maximum of $(1 + e)(GM/h)$ and a minimum of $(1 - e)(GM/h)$.
- When $e = 1$, the origin of velocity space is at the circumference of the hodograph and the magnitude of the velocity vanishes at this point. In this case, the particle goes from a finite distance of closest approach to infinity, reaching infinity with zero speed. Clearly, $e = 1$ implies

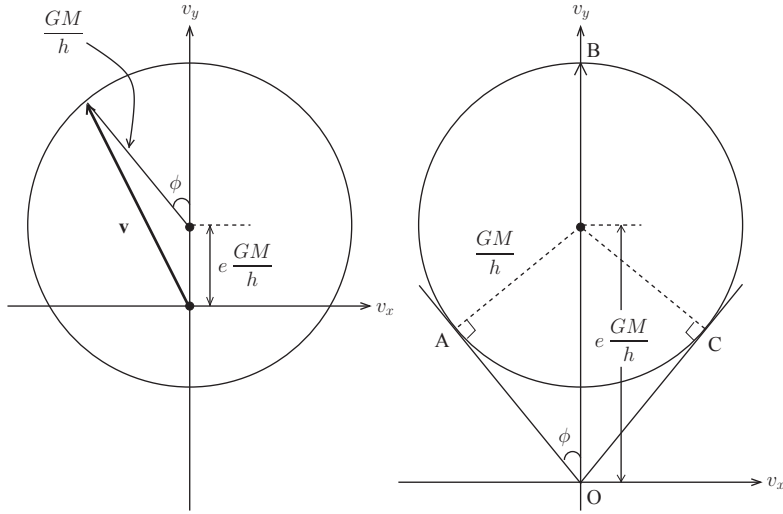


Fig. 3.1: Left: The orbit of a planet in the velocity space moving under the action of an attractive ($1/r^2$) force. This is a circle with center at $(0, e(GM/h))$ and radius GM/h . Here e is just a constant, M is the mass of the Sun and h is the conserved angular momentum per unit mass. Note that the circle is displaced with respect to the origin making the velocity of the planet vary between a maximum and minimum values as long as $e < 1$. This figure is drawn for $e < 1$. Right: Orbit in the velocity space, as in the left figure, but for the case of $e > 1$. Only part of the circular arc is relevant for the motion of the planet which is now moving in an unbounded trajectory in real space.

$u^2 = 2GM/r_{\text{initial}}$ which is just the text book condition for escape velocity.

- When $e > 1$, the origin of velocity space is outside the hodograph and Fig. 3.1 shows the behaviour in this case. The maximum velocity achieved by the particle is OB when the particle is at the point of closest approach in real space. The asymptotic velocities of the particle are OA and OC obtained by drawing the tangents from O to the circle. From the figure it is clear that $\sin \phi_{\text{max}} = e^{-1}$. During the *unbound motion* of the particle, the velocity vector traverses the part ABC . It circles all the way! (Incidentally, the minor arc AC of the hodograph represents the motion under repulsive inverse square force; using the geometrical tricks of this chapter, you should be able to obtain the Rutherford scattering formula from the hodograph.)

Try it out!

An intuitive way of understanding why the hodograph is a circle is as follows: Let us divide the total angle 2π into N equal parts with $\delta\phi \equiv (2\pi/N)$ with a very large N . Let the position of the particle be \mathbf{r}_n when the angle is $\phi_n \equiv n\delta\phi = n(2\pi/N)$. In this discretised version, the particle moves by an amount $\delta\phi$ jumping from one vertex of a large polygon to another in real space in a time interval δt . The corresponding jump

The geometrical reason

in the velocity is by $\delta \mathbf{v}_n = -(GM/r^2)\mathbf{e}_{r_n}\delta t$ according to Newton's law, while $\delta\phi = (h/r^2)\delta t$ from the conservation of angular momentum expressed as $r^2\dot{\phi} = h$. So, $\delta \mathbf{v}_n = (GM/h)\mathbf{e}_{r_n}\delta\phi$. Clearly the *magnitude* of the change in the velocity is a constant equal to $(GM/h)(2\pi/N)$; and direction changes always by the same angle because $\delta \mathbf{v}_n \cdot \delta \mathbf{v}_{n+1} \propto \mathbf{e}_{r_n} \cdot \mathbf{e}_{r_{n+1}} = \delta\phi$. Therefore, it is clear that, $\delta \mathbf{v}_n$ takes the velocity vector from one vertex of a *regular* polygon to next vertex in the velocity space. In the limit of $N \rightarrow \infty$, $\delta\phi \rightarrow 0$ the polygon becomes a circle.

All these must have convinced you that there is something magical about the inverse-square force which is worth exploring. One nice way of understanding the peculiar features of the Kepler (or Coulomb; we will use these interchangeably) problem is to start with a slightly more general potential — which does *not* have these peculiar features — and treat the Kepler problem as a special case of this more general situation. This can be done in many different ways and I will choose to study the dynamics under the action of the potential given by

$$U(r) = -\frac{\alpha}{r} + \frac{\beta}{r^2}, \quad (3.6)$$

which, of course, reduces to the attractive Coulomb/Kepler potential when $\beta \rightarrow 0^+$. For the sake of definiteness, I will take $\alpha > 0$ and $\beta \geq 0$ though most of the analysis can be generalized to other cases.

The classical motion of a particle of mass m , in 3-dimensions, under the action of $U(r)$ is straightforward to analyze using the standard textbook description of a central force problem. Just for fun, let us do it in a slightly different manner. We know that, as with any central force problem, angular momentum \mathbf{J} is conserved, confining the motion to a plane which we will take to be $\theta = \pi/2$. Using $J = mr^2\dot{\phi}$, the energy of the particle can be expressed as

$$E = \frac{1}{2}m\left(\dot{r}^2 + \frac{J^2}{m^2r^2}\right) - \frac{\alpha}{r} + \frac{\beta}{r^2}. \quad (3.7)$$

Combining the two terms with $(1/r^2)$ dependence into C^2/r^2 where $C^2 = (J^2/2m) + \beta$ and completing the square, we get the relation

$$E + \frac{\alpha^2}{4C^2} = \frac{1}{2}m\dot{r}^2 + \left(\frac{C}{r} - \frac{\alpha}{2C}\right)^2 \equiv \mathcal{E}^2, \quad (3.8)$$

where \mathcal{E} is another constant. This suggests introducing a function $f(t)$ via the equations

$$\sqrt{\frac{m}{2}}\dot{r} = \mathcal{E} \sin f(t); \quad \left(\frac{C}{r} - \frac{\alpha}{2C}\right) = \mathcal{E} \cos f(t). \quad (3.9)$$

*Kepler as a limit
of non-Kepler*

*Solvable at no extra
cost since there
is always a J^2/r^2
term!*

Differentiating the second equation with respect to time and using the first equation will give you an expression for \dot{f} . Dividing this expression by $\dot{\phi} = J/mr^2$ leads to the simple relation $(df/d\phi) = (2mC^2/J^2)^{1/2}$. Therefore, f is a linear function of ϕ and from the second equation in Eq. (3.9) we get the equation to the trajectory to be

$$\frac{(2C^2/\alpha)}{r} = 1 + \left(\frac{2\mathcal{E}C}{\alpha} \right) \cos(\omega\phi), \quad (3.10)$$

where

$$\omega^2 = \frac{2m}{J^2} C^2 = \left(1 + \frac{2m\beta}{J^2} \right). \quad (3.11)$$

More generally, we get $(\phi - \phi_0)$ instead of ϕ in Eq. (3.10); we have oriented the axes to set $\phi_0 = 0$.

Now that we have solved the problem completely, let us look at the properties of the solution. To begin with, let us ask what kind of orbit we would expect given the known symmetries of the problem. A particle moving in 3 space dimensions has a phase space which is 6 dimensional. For any time independent central force, we have constancy of energy E and angular momentum \mathbf{J} . Conservation of these four quantities (E, J_x, J_y, J_z) confines the motion to a region of $6 - 4 = 2$ dimensions. The projection of this phase space trajectory on to the xy plane will, in general, fill a two dimensional region of space. So you would expect the orbit to fill a finite two dimensional region of this plane, if there are no other conserved quantities. This is precisely what we find from Eq. (3.10) for a generic value of the conserved quantities J and E . Because ω will not be an integer, when ϕ changes by 2π , the cosine factor will pick up a term $\cos(2\pi\omega)$ which will not be unity. In general, the orbit will fill a 2-dimensional region in the plane between two radii r_1 and r_2 . (See Fig. 3.2.)

Let us count

We can now see how the Kepler (Coulomb) problem becomes rather special. In this case, we have $\beta = 0$ making $\omega = 1$. The curve in Eq. (3.10) closes on itself for *any* value of J and E and — in fact — becomes an ellipse with the latus-rectum $p = (2C^2/\alpha)$ and eccentricity $e = (2\mathcal{E}C/\alpha)$. (You can verify that this is indeed the standard textbook solution to the Kepler problem.) So when $\beta = 0, \omega = 1$ the orbit closes and becomes a one-dimensional curve rather than filling a 2-dimensional region. This analysis shows how turning on a non-zero β completely changes the topological character of the orbit.

The collapse of a dimension

In the argument given above, we linked the nature of the orbit to the number of conserved quantities for the motion. Given the fact that $\beta = 0$ reduces the dimension of the orbital space by one, we will expect to have one more conserved quantity in the problem when $\beta = 0$ which does not exist for $\beta \neq 0$. But we already know one such extra constant which exists for $\beta = 0$ and not otherwise! This is precisely $\mathbf{w} = \mathbf{v} - (GM/h)\mathbf{e}_\phi$

The \mathbf{w} comes to the rescue!

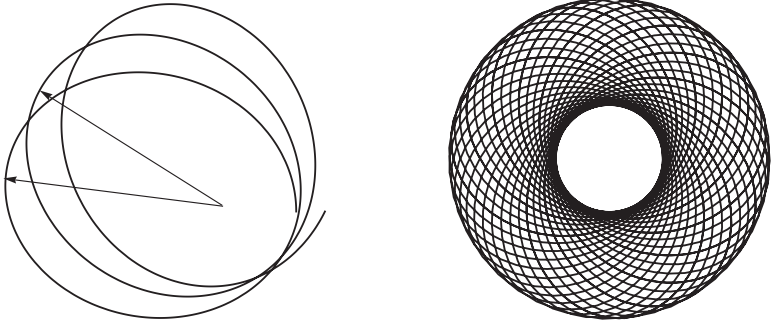


Fig. 3.2: (a) The precessing ellipse as a solution to Eq. (3.10). The vector shows the direction of major axis which precesses. (b) Over a span of time, the orbit fills a 2-dimensional region in the plane for generic values of the parameters. The hodograph in velocity space shows similar behaviour with the velocity vector filling an annular region.

which we discovered in Eq. (3.2). But we needed only one constant of motion while we now have got 3 components of \mathbf{w} which will prevent the particle from moving at all! Such an overkill is avoided because \mathbf{w} satisfies the two, easily verified, constraints because of which it has only one independent component. First, it is obvious that $\mathbf{w} \cdot \mathbf{J} = 0$ because \mathbf{w} is in the orbital plane; this reduces the number of independent components of \mathbf{w} from 3 to 2. Second, its magnitude w can be expressed in terms of E and h and thus is not an independent constant. This is easily seen as follows: Writing $(1/r) = \mathbf{v} \cdot \mathbf{e}_\phi / h$ (which is a cute trick), the conserved energy of the particle is given by

$$\begin{aligned} E &= \frac{1}{2}m \left(v^2 - \frac{2GM}{h} \mathbf{v} \cdot \mathbf{e}_\phi \right) = \frac{1}{2}m \left([\mathbf{v} - (GM/h)\mathbf{e}_\phi]^2 - \frac{G^2M^2}{h^2} \right) \\ &= \frac{1}{2}m \left(w^2 - \frac{G^2M^2}{h^2} \right), \end{aligned} \quad (3.12)$$

showing w^2 is a constant given by

$$w^2 = \frac{2E}{m} + \frac{G^2M^2}{h^2}. \quad (3.13)$$

This further reduces the number of independent constants constrained in \mathbf{w} from 2 to 1, exactly what we needed. It is this extra constant that keeps the planet on a closed orbit. The natural question which arises at this stage is the following: What does this constant mean, geometrically or physically? We will now discuss this issue.

We have developed the entire theory rather trivially by using the natural constant \mathbf{w} . In standard literature, one often uses another constant \mathbf{A}

closely related to \mathbf{w} . To motivate this constant, we consider the following question: Is it possible to represent both the hodographic circle (which lives in the *velocity space*) and the orbital ellipse (which lives in the *real space*) together in the real space itself? To do this sensibly we need to address two issues:

*Shadow of
hodograph in
real space*

(i) Figure 3.1 shows that, in the *velocity space*, the angle ϕ is measured from the v_y axis while in the *real space* the angle ϕ is measured from the x -axis. This tells you that, if we want to plot the hodographic circle as well as the orbital ellipse in the same space, it is more natural to use a vector rotated by 90 degrees with respect to the velocity vector. This can be easily done by taking the cross product of the velocity vector with a unit vector in the direction of \mathbf{J} .

(ii) The vectors \mathbf{v} and \mathbf{r} , of course, have different dimensions and we need to take care of it. This can be done by multiplying the velocity vector by $J/|E|$.

These two facts together suggests defining and studying a vector $(\mathbf{J} \times \mathbf{w})/|E|$ rather than \mathbf{w} as a conserved vector. A simple calculation shows that

$$\begin{aligned} \mathbf{R}_F &\equiv \frac{\mathbf{J} \times \mathbf{w}}{|E|} = \frac{\mathbf{J} \times \mathbf{v}}{|E|} - \frac{GMm}{|E|} (\hat{\mathbf{z}} \times \mathbf{e}_\phi) = \frac{1}{|E|} (\mathbf{J} \times \mathbf{v} + GMm \mathbf{e}_r) \\ &= -\frac{1}{m|E|} (\mathbf{p} \times \mathbf{J} - GMm^2 \mathbf{e}_r) \equiv -\frac{\mathbf{A}}{m|E|} = \frac{\mathbf{A}}{mE}. \end{aligned} \quad (3.14)$$

To arrive at the second equality, we have used Eq. (3.2) and $\mathbf{J} = J\hat{\mathbf{z}}$; to obtain the third equality, we have used the fact $(\hat{\mathbf{z}} \times \mathbf{e}_\phi) = -\mathbf{e}_r$. The fifth equality defines the vector \mathbf{A} (called Runge-Lenz vector).

*... though discovered
by several others;
see Box 3.1*

The conventional route to Runge-Lenz vector starts by computing the time derivative of the quantity $(\mathbf{p} \times \mathbf{J})$ in any central force $f(r)\mathbf{e}_r$ and obtaining:

$$\frac{d}{dt}(\mathbf{p} \times \mathbf{J}) = -mf(r)r^2 \frac{d\mathbf{e}_r}{dt}. \quad (3.15)$$

The miracle of inverse square force is now again in sight: When $f(r)r^2 = \text{constant} = -\alpha$, (with $\alpha = GMm$ in our case) we find that the vector:

$$\mathbf{A} \equiv \mathbf{p} \times \mathbf{J} - \alpha m \mathbf{e}_r \quad (3.16)$$

is conserved. For future reference, let us note the two easily derived properties of \mathbf{A} .

$$A^2 = 2mJ^2E + \alpha^2m^2; \quad \mathbf{A} \cdot \mathbf{J} = 0, \quad (3.17)$$

which are equivalent to Eq. (3.13) and $\mathbf{w} \cdot \mathbf{J} = 0$.

The vector \mathbf{R}_F has direct physical interpretation unlike \mathbf{A} (which, alas, is what people tend to use). The \mathbf{R}_F points from the center of attraction (which is at one focus of the ellipse) to the other focus! It is easy to see

(for example, using the initial conditions) that \mathbf{R}_F points along the major axis away from the center of attraction. Its magnitude is

$$\frac{A^2}{m^2|E|^2} = -\frac{2J^2}{m|E|} + \frac{G^2M^2m^2}{|E|^2} = \left(\frac{GMm}{|E|}\right)^2 \left[1 - \frac{2|E|J^2}{G^2M^2m^3}\right]. \quad (3.18)$$

The first factor in the final expression is $4a^2$ where a is the semi-major axis and the second factor is e^2 . Therefore, $R_F = 2ae$ which is precisely the distance between two foci of the ellipse.

Using \mathbf{R}_F we can write the velocity of the particle in the form

$$\frac{\mathbf{v} \times \mathbf{J}}{|E|} = -\mathbf{R}_F + R_{\max} \mathbf{e}_r; \quad R_{\max} \equiv \frac{GMm}{|E|} = 2a. \quad (3.19)$$

The second relation defines R_{\max} which is the maximum distance from the center of attraction which a particle of energy $-|E|$ can wander off. Eq. (3.19) a remarkable relation in more than one way. To begin with, it gives you the velocity (re-scaled by $J/|E|$ and rotated by 90 degrees, for reasons we explained earlier) in terms of vectors in coordinate space. Second, it shows that the hodograph, when brought back to real space, has a simple geometrical relationship to the elliptical orbit — which is shown in Fig. 3.3. We have now achieved our ambition of drawing both the elliptical orbit and the velocity orbit in the same space. In the process, we have discovered a nice vector \mathbf{R}_F proportional to the Runge-Lenz vector used in the literature.

Incidentally, Fig. 3.3 can be used to provide an elegant (and purely geometrical) derivation for the fact that the orbit is an ellipse. Let us start with the center of force F_1 and draw the position vector \mathbf{r} and the velocity vector \mathbf{v} of the particle at some time t . We are further given that \mathbf{v} satisfies the first relation in Eq. (3.19) which is the same as the result in Eq. (3.2); so we have already used the fact that the particle is moving in an inverse square law force. We will now draw a circle centered at F_1 with the radius $R_{\max} = (GMm/|E|)$ where $-|E|$ is the conserved energy of the particle. We project the point P on to this circle getting P' . At this stage, we have fixed (in the Fig. 3.3) F_1, P, P' and a vector \mathbf{v} representing the velocity of the particle.

We now reflect the vector PP' (denoted, say, by $\ell \equiv R_{\max} \mathbf{e}_r - \mathbf{r}$) on the velocity vector \mathbf{v} . In general, reflecting a vector \mathbf{q} on a plane with unit normal $\hat{\mathbf{n}}$ leads to the vector $\tilde{\mathbf{q}} = \mathbf{q} - 2(\mathbf{q} \cdot \hat{\mathbf{n}})\hat{\mathbf{n}}$. In our case, the relevant normal can be taken to be

$$\hat{\mathbf{n}} = \hat{\mathbf{z}} \times (\mathbf{v}/v) = \frac{1}{Jv}(\mathbf{J} \times \mathbf{v}) = \frac{|E|}{Jv}(\mathbf{R}_F - R_{\max} \mathbf{e}_r). \quad (3.20)$$

so that we get $\tilde{\ell} = \mathbf{R}_F - \mathbf{r}$. That is,

$$\mathbf{r} + \tilde{\ell} = \mathbf{R}_F \quad (3.23)$$

is a constant vector! From the Fig. 3.3 we see that, as P moves along its — as yet unknown — trajectory and P' moves along the circle, F_2 remains unchanged. Since, by construction, $PF_2 = PP'$, it follows that the sum of the lengths $F_1P + F_2P$ remains constant as P moves. This is precisely an ellipse with foci at F_1 and F_2 .

Box 3.1: Runge-Lenz vector, or is it?

Any conserved quantity is of considerable significance in physics and this is especially the case in a problem as important as the one of planetary motion. Given the fact that the existence of the conserved vector \mathbf{A} immediately solves the Kepler problem, it would be interesting to ask who discovered it. Its history is quite fascinating.

It appears that the magnitude of this vector, as a conserved quantity, first appeared in the work of Jacob Hermann [4] in the year 1710. His work was generalized by Bernoulli [5] in the same year, making it a vector, by giving it a direction and magnitude. By the end of the century, Laplace [6] rediscovered the conservation of \mathbf{A} working everything out analytically rather than geometrically. The next important contribution was from Hamilton [7] who, in the middle of nineteenth century, derived this vector in a slightly modified form so that its magnitude is equal to the eccentricity of the orbit. He called it the eccentricity vector and also used it to obtain the hodograph for the Kepler motion.

This vector appears later on in two vector analysis text books, one by Gibbs [8] and another one by Carle Runge [9]. Years later, in 1924, Wilhelm Lenz [10] used this vector for a quantum mechanical treatment of the hydrogen atom. (We will discuss a modern version of this in Chap. 4.) Runge makes no claim of originality in his text book but Lenz refers to Runge in his work. Later on, when Pauli was using a similar technique for the hydrogen atom, he refers to it as “previously utilized by Lenz” and from then on the name Runge-Lenz vector stuck though both Runge and Lenz are quite innocent of *discovering* this vector!

Given the importance to planetary motion, it is also good to know who did *not* discover it: Newton did not, in spite of all his geometrical insights and analytical ingenuity! Scholars have looked in vain for something like the vector \mathbf{A} in Newton’s works hoping that he might have recognized it but that does not seem to be the case.

*Name one person
who did not discover
it!*

The way we approached the problem also shows that the second focus F_2 of the ellipse is not without some significance, as is sometimes thought to be. The Fig. 3.3 shows that the vector F_2P' contains the information about the velocity! Incidentally, using the vector relations obtained above, one can also compute the rate dA_2/dt at which area is swept by the vector F_2P as the planet moves on the elliptical orbit. It is easy to show that the area swept out is proportional to the Jacobi-Mapertuis action (see Eq. (2.39)) for the system:

$$A_2 = \frac{J}{4|E|} \int v^2 dt = \frac{J}{4|E|} \int v d\ell. \quad (3.24)$$

Easy when you note that A_2 is the sum of area of ΔF_1F_2P and the standard area swept by the planet

The second focus F_2 plays a role in this too. (This result has appeared in some classic textbooks with a fairly complicated derivation; see, for example, Ref. [11]. The approach described above provides a simple way of obtaining this result.)

Usually, one associates the conservation laws with the existence of certain symmetries. We know that the time translation symmetry of the Lagrangian leads to energy conservation, spatial translation leads to momentum conservation and rotational invariance leads to the conservation of angular momentum. More generally, consider the variation $\mathbf{x}(t) \rightarrow \mathbf{x}(t) + \delta\mathbf{x}(t)$. The corresponding change in the velocity is given by $\mathbf{v}(t) \rightarrow \mathbf{v}(t) + \delta\mathbf{v}(t)$ where $\delta\mathbf{v}(t) = d\delta\mathbf{x}(t)/dt$. Suppose you can find a particular $\delta\mathbf{x}(t)$ (which is a function of \mathbf{x}, \mathbf{v}) so that, under such a variation, the Lagrangian changes only by a total derivative; that is, $\delta L = dF/dt$ where F is a function of (\mathbf{x}, \mathbf{v}) . This result should arise purely from the structure of the original Lagrangian *without using equations of motion*. On the other hand, explicit variation of the Lagrangian gives

Very simple, but quite useful!

$$\delta L = (\mathbf{f} - \dot{\mathbf{p}}) \cdot \delta\mathbf{x} + \frac{d(\mathbf{p} \cdot \delta\mathbf{x})}{dt}; \quad \mathbf{f} \equiv \frac{\partial L}{\partial \mathbf{x}}; \quad \mathbf{p} \equiv \frac{\partial L}{\partial \dot{\mathbf{x}}}. \quad (3.25)$$

Equating this to dF/dt , we get

$$\frac{dC}{dt} = (\dot{\mathbf{p}} - \mathbf{f}) \cdot \delta\mathbf{x}; \quad C \equiv \mathbf{p} \cdot \delta\mathbf{x} - F. \quad (3.26)$$

So we find that — when the equations of motion $\dot{\mathbf{p}} = \mathbf{f}$ hold — we have the conservation law for the quantity

$$C = \mathbf{p} \cdot \delta\mathbf{x} - F. \quad (3.27)$$

The difficult part, of course, is to find a $\delta\mathbf{x}$ such that the right hand side of Eq. (3.26) is indeed a time derivative.

As a rather trivial example, consider $\delta\mathbf{x} = \boldsymbol{\Omega} \times \mathbf{x}$ which would represent rotation of the coordinates about a direction characterized by $\boldsymbol{\Omega}$ which is assumed to be an infinitesimal vector. If the potential is spher-

ically symmetric, \mathbf{f} will be in the direction of \mathbf{x} and hence $\mathbf{f} \cdot \delta\mathbf{x}$ will vanish. On the other hand, we have

$$\dot{\mathbf{p}} \cdot \delta\mathbf{x} = \dot{\mathbf{p}} \cdot (\boldsymbol{\Omega} \times \mathbf{x}) = \boldsymbol{\Omega} \cdot (\mathbf{x} \times \dot{\mathbf{p}}) = \frac{d}{dt}(\boldsymbol{\Omega} \cdot \mathbf{J}), \quad (3.28)$$

where we have used the fact $\dot{\mathbf{J}} = \mathbf{x} \times \dot{\mathbf{p}}$ and that $\boldsymbol{\Omega}$ is a constant vector. Equation (3.26) now tells you that $C = \boldsymbol{\Omega} \cdot \mathbf{J}$ is conserved for all $\boldsymbol{\Omega}$ which, in turn, requires \mathbf{J} to be conserved. This is the well known result that rotational invariance lead to conservation of angular momentum.

Is there a symmetry of the Lagrangian corresponding to some $\delta\mathbf{x}$ which will lead to the conservation of Runge-Lenz vector? Indeed there is, though this is a bit of a non-trivial transformation given by:

$$\delta\mathbf{x} = \mathbf{a} \times (\mathbf{p} \times \mathbf{x}) + \mathbf{x} \times (\mathbf{p} \times \mathbf{a}), \quad (3.29)$$

where \mathbf{a} is an arbitrary, constant, infinitesimal vector (like the $\boldsymbol{\Omega}$ in the previous example). To discover this transformation, we can use Eq. (3.26) to do a bit of reverse-engineering! Suppose you do know that a specific function $C(\mathbf{x}, \mathbf{p})$ is indeed conserved. Then if you compute dC/dt and group together the terms in the form of Eq. (3.26), you can read off $\delta\mathbf{x}$. In the case Runge-Lenz vector, we take $C = \mathbf{a} \cdot \mathbf{A}$ and compute dC/dt , using the identity

$$\frac{d\mathbf{e}_r}{dt} = -\frac{1}{r^3}((\dot{\mathbf{x}} \times \mathbf{x}) \times \mathbf{x}) = -\frac{1}{mr^3}(\mathbf{x} \times \mathbf{J}) \quad (3.30)$$

and carefully group together terms involving $\dot{\mathbf{p}}$. This gives

$$\begin{aligned} \mathbf{a} \cdot \dot{\mathbf{A}} &= \mathbf{a} \cdot (\dot{\mathbf{p}} \times \mathbf{J}) + \mathbf{a} \cdot (\mathbf{p} \times (\mathbf{x} \times \dot{\mathbf{p}})) + \frac{GMm}{r^3}(\mathbf{x} \times \mathbf{J}) \cdot \mathbf{a} \\ &= \dot{\mathbf{p}} \cdot (\mathbf{J} \times \mathbf{a}) + (\mathbf{a} \cdot \mathbf{x})(\dot{\mathbf{p}} \cdot \mathbf{p}) - (\dot{\mathbf{p}} \cdot \mathbf{a})(\mathbf{p} \cdot \mathbf{x}) + (-\mathbf{f} \cdot (\mathbf{J} \times \mathbf{a})) \\ &= \dot{\mathbf{p}} \cdot [(\mathbf{J} \times \mathbf{a} + \mathbf{p}(\mathbf{a} \cdot \mathbf{x}) - \mathbf{a}(\mathbf{p} \cdot \mathbf{x})) - \mathbf{f} \cdot (\mathbf{J} \times \mathbf{a})] \\ &= \dot{\mathbf{p}} \cdot [(\mathbf{J} \times \mathbf{a} + \mathbf{x} \times (\mathbf{p} \times \mathbf{a})) - \mathbf{f} \cdot (\mathbf{J} \times \mathbf{a})] \\ &= \dot{\mathbf{p}} \cdot \delta\mathbf{x} - \mathbf{f} \cdot \delta\mathbf{x} \end{aligned} \quad (3.31)$$

with

$$\delta\mathbf{x} = \mathbf{J} \times \mathbf{a} + \mathbf{x} \times (\mathbf{p} \times \mathbf{a}). \quad (3.32)$$

We can now see that $\delta\mathbf{x}$ is indeed given by the expression in Eq. (3.29). Note that the second term in $\delta\mathbf{x}$ is perpendicular to \mathbf{x} and does not contribute to the $\mathbf{f} \cdot \delta\mathbf{x}$ in Eq. (3.31). One can remove the arbitrary vector \mathbf{a} , if one wants and write down an expression for the infinitesimal transformation relevant to the conservation of the s -th component of \mathbf{A} . It is given by

$$\delta x_i = \frac{\varepsilon}{2} [2p_i x_s - x_i p_s - \delta_{is}(\mathbf{x} \cdot \mathbf{p})]; \quad (i = 1, 2, 3), \quad (3.33)$$

Formally, the symmetry transformation δq^i corresponding to a constant of motion $C(q, p)$ is given by the Poisson bracket $\delta q^i = \{q^i, C\}$, but let us keep it simple.

where ε is a constant infinitesimal parameter. That is, you change (x_1, x_2, x_3) in the above manner keeping s fixed at some value 1 or 2 or 3. (It is probably nicer to think in terms of Eq. (3.29).) You may wonder why such a strange $\delta\mathbf{x}$ exists. It is possible to relate this to rotations in a fictitious 4-dimensional space (see Chapter 4) but unfortunately it does not make it any more enlightening.

We can close the logical loop by asking what happens to the eccentricity vector when we add a β/r^2 term to the $1/r$ potential. Obviously, if you add a $1/r^3$ component to the force, (which can arise, for example, from the general relativistic corrections or because the Sun is not spherical and has a small quadrupole moment) \mathbf{J} and E are still conserved but not \mathbf{A} . If the perturbation is small, it will make the direction of \mathbf{A} slowly change in space and we will get a “precessing” ellipse, which will of course fill a 2-dimensional region. For the potential in Eq. (3.6) we find, using Eq. (3.15), that the rate of change of Runge-Lenz vector is now given by $\dot{\mathbf{A}} = -(2\beta m/r)(d/dt)(\mathbf{r}/r)$. The change $\Delta\mathbf{A}$ per orbit is obtained by integrating $\dot{\mathbf{A}}dt$ over the range $(0, T)$ where T is the period of the original orbit. Doing one integration by parts and changing the variable of integration from t to the polar angle ϕ , we get $\Delta\mathbf{A}$ per orbit to be

*Change in \mathbf{A}
measures the
precession rate*

$$\Delta\mathbf{A}|_{\text{orbit}} = -2\beta m \int_0^{2\pi} \frac{\mathbf{r}}{r^3} \frac{dr}{d\phi} d\phi. \quad (3.34)$$

Let us take the coordinate system such that the unperturbed orbit originally had \mathbf{A} pointing along the x -axis. After one orbit, a ΔA_y component will be generated and the major axis of the ellipse would have precessed by an amount $\Delta\phi = \Delta A_y/A$. The ΔA_y can be easily obtained from Eq. (3.34) by using $y = r \sin\phi$, converting the dependent variable from r to $u = (1/r)$ and substituting $(du/d\phi) = -(A/J^2) \sin\phi$ [which comes from Eq. (3.3)]. This gives the angle of precession per orbit to be

*Do not confuse
 $u \equiv 1/r$ with
velocity defined
earlier!*

$$\Delta\phi = \frac{\Delta A_y}{A} = \frac{2\beta m}{A} \int_0^{2\pi} \sin\phi \frac{du}{d\phi} d\phi = -\frac{2\pi\beta m}{J^2}. \quad (3.35)$$

Since we have the exact solution in Eq. (3.10), you can easily verify that this is indeed the precession of the orbit when β/r^2 is treated as a perturbation. The Runge-Lenz vector not only allows us to solve the $(1/r)$ problem, but even tells us how an r^{-2} perturbation makes the orbit precess! This is indeed the primary effect when we introduce physically relevant modifications to the Kepler problem.

The first generalization of the Kepler problem that you might think of will be to introduce the effects of special relativity. This turns out to be more non-trivial than one might have imagined for the following reason. In the non-relativistic context, the motion of a particle under the action of a potential V is governed by the equation $dp_\alpha/dt = -\partial_\alpha V$ where $\alpha = 1, 2, 3$ denotes the three spatial components of the momentum \mathbf{p} and ∂_α denotes

*What happens when
relativity steps in?*

Relativity prevents irresponsible addition of interactions

the derivative with respect to the coordinate x^α . One might have thought that the natural generalization of this Newtonian result into special relativistic domain would involve the following replacements: Change the three momentum p_α to the four momentum p_i (with $i = 0, 1, 2, 3$), the coordinate time t into the proper time τ of the particle and the three dimensional gradient ∂_α to the four-gradient ∂_i . This would have led to the equation $dp_i/d\tau = -\partial_i V$. Unfortunately, there is a problem with this “generalization”. The four-velocity u^i satisfies the constraint:

$$u_i u^i = \frac{dx_i dx^i}{d\tau d\tau} = -\frac{d\tau^2}{d\tau^2} = -1 . \quad (3.36)$$

Since the four-momentum is $p_i = mu_i$, we have the constraint

$$u^i \frac{dp_i}{d\tau} = mu^i \frac{du_i}{d\tau} = \frac{m}{2} \frac{d}{d\tau} (u_i u^i) = 0 , \quad (3.37)$$

and we have used Eq. (3.36). This implies that our potential V has to satisfy the constraints $u^i \partial_i V = 0$; that is, the potential should not change along the worldline of the particle which is not possible in general.

This is why forces are velocity dependent in relativistic theories

So the generalization to special relativity has to come from some other direction. One possibility is to note that the “Kepler problem” also arises in electrodynamics when we consider the motion of a test charge in the Coulomb field of another charge. Since we have a fully special relativistic formulation of electrodynamics, we can attempt to study the motion of a particle under the action of a four-vector potential $A^i = (\Phi(r), 0, 0, 0)$ which would correspond to a centrally symmetric electrostatic potential.

Hamilton-Jacobi for central force: general case

The study of orbits in external fields is most economically done using the Hamilton-Jacobi equation which — as we saw in Chapter 2 — has the blessings of quantum theory. Since energy E and angular momentum J will be conserved in all the contexts we consider, the action S can be expressed in the form

$$S(t, r; E, J) = -Et + J\phi + F(r; E, J) , \quad (3.38)$$

where (r, ϕ) are the standard polar coordinates in the plane of orbit and $F(r; E, J)$ has to be determined by integrating the Hamilton-Jacobi equation. The orbital equation $r = r(\phi)$ can be obtained by differentiating S with respect to J and equating it to a constant:

$$\phi + \frac{\partial F}{\partial J} = \phi_0 = \text{constant} . \quad (3.39)$$

The different contexts we would be interested in, differ only in the nature of Hamilton-Jacobi equation; once we obtain the orbital equation in Eq. (3.39) one can compare the different models fairly easily.

We recall that, in the case of standard Newtonian context, for particle moving in a central potential $V(r)$, the Hamilton-Jacobi equation is $\partial S/\partial t + H = 0$. It is easy to show that F satisfies the equation

$$\left(\frac{dF}{dr}\right)^2 = 2m(E - V) - (J^2/r^2). \quad (3.40)$$

This, in turn, allows us to write the orbital equation in Eq. (3.39) in the form

$$\phi - \phi_0 = \int \frac{dr(J/r^2)}{[2m(E - V) - (J^2/r^2)]^{1/2}}. \quad (3.41)$$

Converting this into an equation for $dr/d\phi$ and introducing the standard substitution $u \equiv (1/r)$ we can obtain the differential equation satisfied by $u(\phi)$:

$$u'' + u = -\frac{m}{J^2} \frac{dV}{du}, \quad (3.42)$$

where the prime denotes differentiation with respect to ϕ . In the standard Kepler problem, $V = -GMm/r = -GMmu$ so that the right hand side of Eq. (3.42) becomes a constant and we get the solution $u = \alpha + \beta \cos \phi$ which represents a conic section.

The Newtonian case

For the relativistic particle with charge q moving in an electromagnetic field with $A^i = (\Phi, 0, 0, 0)$ the Hamilton-Jacobi equation is given by Eq. (2.15) and the corresponding differential equation for F given by

$$\left(\frac{dF}{dr}\right)^2 = \frac{1}{c^2}(V - E)^2 - \frac{J^2}{r^2} - m^2 c^2; \quad V(r) \equiv q\Phi(r). \quad (3.43)$$

It is fairly straightforward to show that, in this case, Eq. (3.42) gets modified to the form

Add special relativity

$$u'' + u = -\frac{(E - V)}{J^2 c^2} \left(\frac{dV}{du}\right) = -\frac{E/c^2}{J^2} \frac{dV}{du} + \frac{1}{2} \frac{1}{J^2 c^2} \frac{dV^2}{du}. \quad (3.44)$$

Comparing Eq. (3.44) with Eq. (3.42) we see that the first term involves replacement of m by E/c^2 which, of course, makes sense in relativity; the second term shows that the potential picks up a V^2 term as a correction which can be traced back to the fact that while the square of the momentum is proportional to energy in the non-relativistic case, it is proportional to the square of the energy in special relativity.

More formally, we can attempt to define a Newtonian effective potential V_{eff} using which we will obtain the same equation of motion. In the case of motion in a Coulomb field with $V(r) = -\alpha/r = -\alpha u$ where $\alpha = Qq$, say, this requires us to satisfy the condition

$$\frac{m}{J^2} \frac{dV_{\text{eff}}}{du} = -\frac{\alpha E}{J^2 c^2} - \frac{\alpha^2}{J^2 c^2} u, \quad (3.45)$$

which integrates to give

$$V_{\text{eff}} = - \left(\frac{E}{mc^2} \right) \alpha u - \frac{\alpha^2}{mc^2} \frac{u^2}{2} . \quad (3.46)$$

Since E/mc^2 is γ , we can think of the first term as the original Coulomb potential transformed to the rest frame of the moving body. The second term is a purely relativistic correction. (Of course, V_{eff} is not a ‘genuine’ potential because it depends on the parameters of the particle, like E .) In this case, the orbit equation becomes:

$$u'' + \omega^2 u = \frac{\alpha E}{J^2 c^2}; \quad \omega^2 \equiv 1 - \frac{\alpha^2}{J^2 c^2} . \quad (3.47)$$

The introduction of c by special relativity has led to a new dimensionless combination (α/Jc) . Obviously, we will expect new features — with no non-relativistic analogue — to arise when $(\alpha/Jc) \gtrsim 1$, because ω will be imaginary for $(\alpha/Jc) > 1$. This is indeed true but let us first consider the case of $(\alpha/Jc) < 1$. In this case, the trajectory obtained by solving Eq. (3.47) can be expressed in the form (compare with Eq. (3.10))

*The special
relativistic
trajectory*

$$\frac{1}{r} = \frac{1}{R} \cos(\omega\phi) + \frac{E\alpha}{c^2 J^2 \omega^2} , \quad (3.48)$$

where

$$R \equiv \frac{J\omega^2}{mc} \left[\left(\frac{E}{mc^2} \right)^2 - 1 + \frac{\alpha^2}{c^2 J^2} \right]^{-1/2} , \quad (3.49)$$

is a constant. In a more familiar form, the trajectory is $l/r = (1 + e \cos \omega\phi)$ with

$$l = \frac{c^2 J^2 \omega^2}{E|\alpha|}; \quad e^2 = \frac{J^2 c^2}{\alpha^2} \left[1 - \frac{m^2 c^4 \omega^2}{E^2} \right] . \quad (3.50)$$

It is easy to verify that, when $c \rightarrow \infty$, this reduces to the standard equation for an ellipse in the Kepler problem. In terms of the non-relativistic energy $E_{\text{nr}} \equiv E - mc^2$, we get, to leading order, $\omega \approx 1$, $l \approx J^2/m|\alpha|$ and $e^2 \approx 1 + (2E_{\text{NR}}J^2/m\alpha^2)$ which are the standard results.

*The precession,
again!*

In the fully relativistic case all these expressions change but the key new effect arises from the fact that $\omega \neq 1$. Due to this factor, the trajectory is not closed and the ellipse precesses. (See Fig. 3.2.) When $\omega \neq 1$ the r in Eq. (3.48) does not return to the value at $\phi = 0$ when $\phi = 2\pi$; instead, we need a further turn by $\Delta\phi$ (the ‘angle of precession’) for r to return to the original value. This is determined by the condition $(2\pi + \Delta\phi)\omega = 2\pi$. From Eq. (3.48) we find that the orbit precesses by the angle

$$\Delta\phi = 2\pi \left[\left(1 - \frac{\alpha^2}{c^2 J^2} \right)^{-1/2} - 1 \right] \simeq \frac{\pi\alpha^2}{c^2 J^2} \quad (3.51)$$

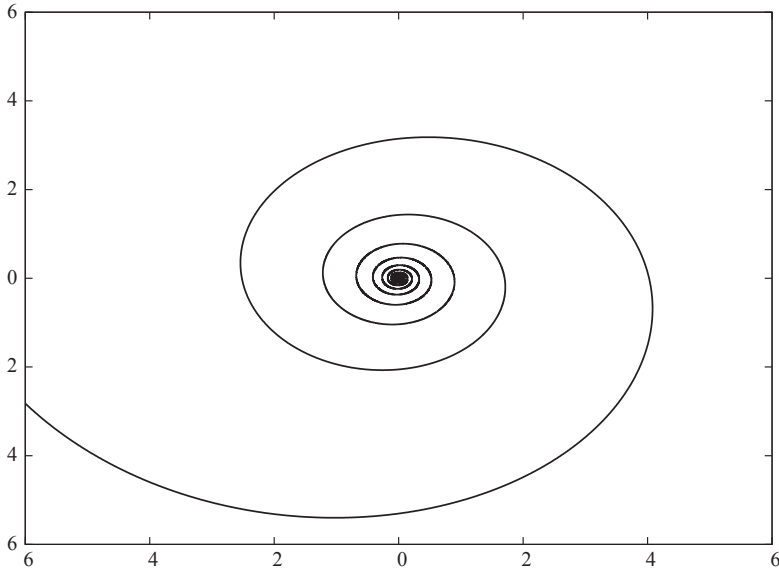


Fig. 3.4: Trajectory of a charged particle around another charge of opposite sign in special relativistic “Kepler problem”. For sufficiently low angular momentum, the trajectory spirals down to the center of attraction. This phenomenon has no non-relativistic analogue.

per orbit where the second expression is valid for $\alpha^2 \ll c^2 J^2$. This is a purely relativistic effect and vanishes when $c \rightarrow \infty$.

There is another peculiar feature which arises in the special relativistic case which has *no* non-relativistic analogue. You would have noticed that ω^2 in Eq. (3.47) has two terms of opposite sign and in obtaining our result, we have tacitly assumed that $\omega^2 > 0$. But in principle, one can have a situation with very low *but non-zero* angular momentum making $\omega^2 < 0$. This is a feature which non-relativistic Kepler problem simply does not have and — under such drastic change of circumstances — one can no longer think in terms of perturbation theory and precessing ellipses. The Eq. (3.47) now has the solution

Here is something new!

$$(\alpha^2 - c^2 J^2) \frac{1}{r} = \pm c \sqrt{(JE)^2 + m^2 c^2 (\alpha^2 - J^2 c^2)} \cdot \cosh \left(\phi \sqrt{\frac{\alpha^2}{c^2 J^2} - 1} \right) - E \alpha. \quad (3.52)$$

Overcoming the angular momentum; a new special relativistic feature!

In this expression, we take the positive root for $\alpha > 0$ and negative root for $\alpha < 0$. It is obvious that, as ϕ increases, $(1/r)$ keeps increasing in the

case of attractive motion so that the test particle spirals to the origin! A typical trajectory is shown in [Figure 3.4](#).

This does not happen for the Kepler problem in Newtonian physics. As is well known, the angular momentum term gives a repulsive J^2/r^2 contribution to the effective potential in the central force problem. In the case of $-(1/r)$ potential, the angular momentum term prevents any particle with non-zero J from reaching the origin. This is not the case in special relativistic motion under attractive Coulomb field. If the angular momentum is less than a critical value, α/c , then the particle spirals down to the origin.

*No, you can't do
gravity this way;
wait till Chapter 11*

If we think of α as GMm , the second term in Eq. (3.46) gives a correction to the potential $(-G^2M^2/2c^2)(m/r^2)$. This term *will* lead to a precession of ellipse but the model is totally wrong. One cannot represent gravity using a vector potential; in such a theory, like charges repel while the gravity has to be attractive. The proper way of generalizing the gravitational Kepler problem, taking into account the effects of relativity, is of course to use general relativity to describe the gravitational field. We will discuss this in Chapter 11.

The Importance of being Inverse-square

4

We saw in Chapter 3 that the motion of a particle in the attractive $(1/r)$ potential has several peculiar features. This potential arises both in the case of planetary motion (the Kepler problem) as well as in the study of atomic systems like the hydrogen atom (the Coulomb problem). In this chapter, we shall complement the classical discussion of Chapter 3 by describing several peculiar *quantum* features [12] that arise in the study of the inverse square law force.

Inverse-square is special in QM as well!

We learnt in Chapter 3 that a simple way to understand the special properties of the inverse square potential is to start with the potential given by

$$U(r) = -\frac{\alpha}{r} + \frac{\beta}{r^2}, \quad (4.1)$$

and study the limit of $\beta \rightarrow 0$. To study quantum mechanics, we first need to solve the Schrödinger equation for the potential in Eq. (4.1). It turns out that this is indeed possible and the analysis proceeds exactly as in the case of normal hydrogen atom problem. Once the angular dependence is separated out using the standard spherical harmonics $Y_{\ell m}(\theta, \phi)$, the radial part of the wavefunction $R(r)$ will satisfy the differential equation

No extra cost; we anyway had an $\ell(\ell+1)/r^2$ in the equation!

$$R'' + \frac{2}{r}R' + \frac{2m}{\hbar^2} \left\{ E - \frac{\hbar^2}{2mr^2} \ell(\ell+1) - \frac{\beta}{r^2} + \frac{\alpha}{r} \right\} R = 0, \quad (4.2)$$

where the prime denotes derivative with respect r and $E(<0)$ is the energy eigenvalue. Let us introduce a new variable ρ by $\rho = 2(-2mE)^{1/2}r/\hbar$ and two new constants s and n by

$$s(s+1) \equiv \frac{2m\beta}{\hbar^2} + \ell(\ell+1); \quad n \equiv \frac{\alpha}{\hbar} \left(\frac{m}{-2E} \right)^{1/2}. \quad (4.3)$$

(The $\ell = 0, 1, 2, \dots$ is the eigenvalue of the angular momentum operator while s is just a parameter; in general, it will *not* be an integer.) The radial

equation can be rewritten as:

$$\frac{d^2 R}{d\rho^2} + \frac{2}{\rho} \frac{dR}{d\rho} + \left(-\frac{1}{4} + \frac{n}{\rho} - \frac{s(s+1)}{\rho^2} \right) R = 0. \quad (4.4)$$

This is *identical* to the standard radial equation for the hydrogen atom which, actually, is to be expected. Algebraically, this arises because the angular momentum term always has a r^{-2} dependence and the β/r^2 part of the potential just combines with the angular momentum term as shown in the first equation in Eq. (4.3).

The quantization condition for energy levels now follows in a straightforward manner (as in the case of usual hydrogen atom) and you will find that $p \equiv (n - s - 1)$ must be a positive integer or zero for well-behaved solutions to exist. (The s is taken to be the positive root of the quadratic equation in Eq. (4.3).) This allows us to obtain the energy levels to be

$$-E = \frac{2\alpha^2 m}{\hbar^2} \left\{ 2p + 1 + \left[(2\ell + 1)^2 + \frac{8m\beta}{\hbar^2} \right]^{1/2} \right\}^{-2}. \quad (4.5)$$

It is now clear that the nature of energy levels depends rather crucially on whether $\beta = 0$ or $\beta \neq 0$. When $\beta \neq 0$ we find that the energy levels depend both on p and ℓ . That is, if we keep p fixed and change ℓ , the energy of the state changes because it depends on both the quantum numbers. On the other hand, when $\beta = 0$, Eq. (4.3) tells us that $s = \ell$. Therefore, $p \equiv (n - s - 1) = (n - \ell - 1)$ and the factor inside the curly bracket in Eq. (4.5) reduces to

$$(2n - 2\ell - 2) + 1 + (2\ell + 1) = 2n. \quad (4.6)$$

In this limit, the energy depends only on the principle quantum number n and becomes independent of the angular quantum number ℓ . The states with same n and different ℓ become degenerate which is the origin of the phrase “accidental degeneracy of the Coulomb potential”. (In a way, this is similar to the classical orbits closing in the case of $\beta = 0$.) As I said before, starting from the potential in Eq. (4.1), solving the problem completely and then taking the limit of $\beta \rightarrow 0$ helps us to distinguish such “accidental” results from more generic results.

In the classical Coulomb problem, we could find the orbit purely algebraically using the Runge-Lenz vector without solving a differential equation. Can we do something similar in the case of quantum mechanics? Can we find the energy levels of the hydrogen atom without explicitly solving the Schrödinger equation? It turns out that this is indeed possible as was first shown by Pauli in 1926. The operator algebra which is involved is straightforward but lengthy and hence I will just indicate the key steps. (One good place to look up the details of the algebra is Ref. [13].)

*Energy levels
depend on ℓ ,
if $\beta \neq 0$...*

*... but independent
of ℓ when $\beta = 0$*

*Runge-Lenz trick
works in QM too!*

We now switch to the $\beta = 0$ (viz., the standard Coulomb problem) and define an operator $\mathbf{M} = \mathbf{A}/m$ corresponding to the classical Runge-Lenz vector (divided by m for convenience). Classically, in the definition of the Runge-Lenz vector, we could have used either $\mathbf{p} \times \mathbf{J}$ or $-\mathbf{J} \times \mathbf{p}$ because $\mathbf{p} \times \mathbf{J} = -\mathbf{J} \times \mathbf{p}$. But this is not true in quantum mechanics because of the non-trivial commutation relations. Hence the appropriate operator — which will be Hermitian — needs to be defined as

*Correct definition
of \mathbf{A} in QM*

$$\mathbf{M} = \frac{1}{2m} (\mathbf{p} \times \mathbf{J} - \mathbf{J} \times \mathbf{p}) - \alpha \frac{\mathbf{r}}{r}, \quad (4.7)$$

where each term is now an operator. By explicit computation, you can verify that the following identities are satisfied:

$$[\mathbf{M}, H] = 0; \quad \mathbf{J} \cdot \mathbf{M} = \mathbf{M} \cdot \mathbf{J} \quad (4.8)$$

and

$$M^2 = \alpha^2 + \frac{2H}{m} (\hbar^2 + J^2), \quad (4.9)$$

where H is the Hamiltonian. You can now recognize the correspondence between the operator relation in Eq. (4.8) and the classical properties of the Runge-Lenz vector given by Eq. (3.17). The relation in Eq. (4.9), however, is a bit non-trivial because it has an additional \hbar^2 term which is purely quantum mechanical and arises because of the non-commuting nature of the operators. Further, we have three commutation rules which can all be directly obtained from the definitions:

*Quantum analogues
of classical results,
almost*

$$\begin{aligned} [J_i, J_j] &= i\hbar \epsilon_{ijk} J_k; \\ [M_i, J_j] &= i\hbar \epsilon_{ijk} M_k; \\ [M_i, M_j] &= -2i(\hbar/m) H \epsilon_{ijk} J_k. \end{aligned} \quad (4.10)$$

The first one is well-known, of course; the second reflects the fact that the components of \mathbf{M} behave as a vector under spatial rotations. The really non-trivial one is the third commutation rule which — by a series of manipulations — allows us to deduce the eigenvalues of H . I will now outline this procedure.

We first note that, since $H, \mathbf{M}, \mathbf{J}$ are conserved quantities (in the sense that they all commute with the Hamiltonian), we can confine ourselves to a sub-space of a Hilbert space that corresponds to a particular eigenvalue $E (< 0)$ of the Hamiltonian H . Working in this subspace, we can replace H by its eigenvalue in the third commutation relation in Eq. (4.10). We then rescale \mathbf{M} by $\mathbf{M}' \equiv (-m/2E)^{1/2} \mathbf{M}$ so that the last two commutation relations in Eq. (4.10) can be expressed in the form

$$[M'_i, J_j] = i\hbar \epsilon_{ijk} M'_k; \quad [M'_i, M'_j] = i\hbar \epsilon_{ijk} J_k, \quad (4.11)$$

A trick worth learning

showing that they constitute a closed set. This set can be separated by defining two other operators $\mathbf{I} = (1/2)(\mathbf{J} + \mathbf{M}')$, $\mathbf{K} = (1/2)(\mathbf{J} - \mathbf{M}')$ which will satisfy the commutation relations:

$$[I_i, I_j] = i\hbar \epsilon_{ijk} I_k; \quad [K_i, K_j] = i\hbar \epsilon_{ijk} K_k, \quad (4.12)$$

We now have two sets of angular momentum operators!

with other commutators vanishing. From our knowledge of the angular momentum operators, we know that the spectra of I^2 and K^2 are given by $j(j+1)\hbar^2, k(k+1)\hbar^2$ where $(j, k) = 0, 1/2, 1, \dots$. But since I and K obey the additional constraints:

$$I^2 - K^2 = \mathbf{J} \cdot \mathbf{M} = 0, \quad (4.13)$$

we only need to consider the subspace with $j = k$. Then the operator

$$\frac{1}{2}(J^2 + M'^2) = \frac{1}{2}[(\mathbf{I} + \mathbf{K})^2 + (\mathbf{I} - \mathbf{K})^2] = I^2 + K^2, \quad (4.14)$$

will have the eigenvalues $[j(j+1) + k(k+1)]\hbar^2 = 2k(k+1)\hbar^2$ (because $j = k$) with $k = 0, 1/2, 1, \dots$. On the other hand, we also have

$$\frac{1}{2}(J^2 + M'^2) = \frac{1}{2}\left[J^2 - \frac{m}{2E}M^2\right] = -\frac{m\alpha^2}{4E} - \frac{1}{2}\hbar^2, \quad (4.15)$$

where the last relation arises from Eq. (4.9). Using the eigenvalues of the operator in Eq. (4.14), we see that E is quantized in the form:

$$E = -\frac{m\alpha^2}{2\hbar^2(2k+1)^2}, \quad (4.16)$$

No solving differential equations!

which is the standard result. So, once again, the existence of an extra conserved quantity allows us to solve the problem completely.

The reason it works

The physical meaning of the above steps relies on the commutation relations in Eq. (4.12) and the constraint $I^2 = K^2$. You can think of the commutation relation in Eq. (4.12) as describing rotations in two different planes in a hypothetical 4-dimensional space with coordinates (q_1, q_2, q_3, q_4) . In other words, the hydrogen atom problem seems to exhibit rotational invariance in a hypothetical 4-dimensional space! In fact, the situation is better than that. You can map the Hamiltonian of the 3-dimensional hydrogen atom to that of a 4-dimensional isotropic harmonic oscillator with an extra restriction which comes from the condition $I^2 = K^2$. I will now describe how this comes about. Since this mapping is somewhat complicated mathematically, we will do this in steps.

A cute trick

Let us begin with the Hamilton-Jacobi equation for the central force and consider the radial part of the action which obeys Eq. (3.40). When

$V(r) = Ar^k$, this equation reads as:

$$\left(\frac{dF}{dr}\right)^2 \equiv p_r^2 \equiv 2m \left(E - Ar^k - \frac{J^2}{2mr^2} \right). \quad (4.17)$$

Suppose we now change variables from r to s such that $r = s^n$. Elementary algebra now leads to a modified form of Eq. (4.17) given by

$$\left(\frac{dF}{ds}\right)^2 = 2mn^2 \left(Es^{2n-2} - As^{n(k+2)-2} - \frac{J^2}{2ms^2} \right). \quad (4.18)$$

You now notice that if we rescale J by Jn , the last term in Eq. (4.18) has the same structure as the last term in Eq. (4.17) (with s treated as a radial coordinate) and represents the contribution due to the angular momentum. As regards the other two terms in Eq. (4.18), we would like one of them to be a constant representing the energy, say, \mathcal{E} of the system while the other term should represent some central potential, say, $\mathcal{V}(s)$. When $n = 1$ and $r = s$ — which is the original system — the first term in Eq. (4.18) represents the energy, while the second term represents the potential. But there is another possibility: If we choose $n = 2/(2+k)$, we can make the second term in Eq. (4.18) a constant. For this choice, we will have

$$2n - 2 = \frac{4}{2+k} - 2 = -\frac{2k}{2+k}, \quad (4.19)$$

and the first term in Eq. (4.18) will correspond to a potential which is another power law. In that case, Eq. (4.18) becomes

$$\left(\frac{dF}{ds}\right)^2 = 2m \left(\bar{E}s^{-2k/(k+2)} - \bar{A} - \frac{\bar{J}^2}{2ms^2} \right), \quad (4.20)$$

where $\bar{E} = n^2E$, $\bar{A} = n^2A$, $\bar{J} = nJ$ are rescaled parameters of the problem. This represents the relevant equation for the radial action for another central force problem (in the variable s) with energy \mathcal{E} and potential $\mathcal{V}(s)$ where

$$\mathcal{E} = -\bar{A}; \quad \mathcal{V}(s) = -\bar{E}s^{-2k/(k+2)}. \quad (4.21)$$

Let us specialize now to the Coulomb problem with $k = -1$ and $A = -Zq^2$ where q is the charge of the electron and Z is the atomic number. Let $E = -|E|$ be the negative energy corresponding to a bound state. In this case, $n = 2/(2+k) = 2$ and Eq. (4.19) gives $[-2k/(2+k)] = 2$. We now see from Eq. (4.21) that the original problem gets mapped to another central force problem with

$$\mathcal{E} = -\bar{A} = 4Zq^2; \quad \mathcal{V}(s) = -4Es^2 = 4|E|s^2. \quad (4.22)$$

Voila!

We have transformed the Coulomb problem to a harmonic oscillator! A parameter describing the original potential $4Zq^2$ appears as the energy of the oscillator and the original bound state energy appears as the squared frequency of the oscillator.

*The quantum case
in D dimensions*

The same idea works in quantum theory for the Coulomb problem in $D = 3$, if the oscillator is in $D = 4$. To see this, let us consider an isotropic harmonic oscillator in a hypothetical D -dimensional space with coordinates (q_1, q_2, \dots, q_D) . Let us introduce in this space the standard radial coordinate s with $s^2 = q^i q_i$ and $(D - 1)$ angular coordinates $(\theta_1, \theta_2, \dots, \theta_{D-1})$. (This is just a generalization of what we would have done in $D = 3$ in terms of one radial coordinate r and two angular coordinates θ and ϕ .) The Hamiltonian for a quantum isotropic oscillator will be the sum of kinetic and potential energy terms where the potential energy is just $\mathcal{V}(s) = (1/2)m\Omega^2 s^2$, where m is the mass of the particle and Ω is the angular frequency of the oscillator. The quantum mechanical operator for the kinetic energy part $\hat{\mathbf{p}}^2/2m = -(\hbar^2/2m)\nabla_{(D)}^2$ — where $\nabla_{(D)}^2$ is the Laplacian in D dimensions — can be separated into a radial part involving $\hat{\mathbf{p}}_s^2$ and an angular part having the form $\hat{\mathbf{L}}^2/s^2$ where $\hat{\mathbf{L}}^2$ is the Laplacian on the $(D - 1)$ sphere defined by $s = \text{constant}$. (This is again in complete analogy with what we do in $D = 3$. There, we would have separated the radial and angular parts of the Laplacian ∇^2 in exactly the same way.) The relevant Schrödinger equation will now read as:

$$\left\{ \frac{1}{2m} \left[\hat{\mathbf{p}}_s^2 + \frac{\hat{\mathbf{L}}^2}{s^2} \right] + \frac{1}{2}m\Omega^2 s^2 - E_{\text{osc}} \right\} \Psi = 0, \quad (4.23)$$

where E_{osc} is the energy eigenvalue of the $D = 4$ oscillator.

Let us separate out the angular and radial parts of the wavefunction Ψ as $\Psi(s, \theta_i) = \mathcal{R}(s)\Phi(\theta_i)$ with $\hat{\mathbf{L}}^2\Phi = L^2\Phi$ where L^2 is the relevant eigenvalue of the angular Laplacian. Concentrating on the radial equation, we will play the old trick and introduce the variable $\rho \equiv s^2$ and divide Eq. (4.23) throughout by ρ . This leads to the equation

$$\left\{ \frac{1}{2m} \left[\frac{\hat{\mathbf{p}}_s^2}{\rho} + \frac{L^2}{\rho} \right] - \frac{E_{\text{osc}}}{\rho} + \frac{1}{2}m\Omega^2 \right\} \mathcal{R} = 0. \quad (4.24)$$

If you compare Eq. (4.23) and Eq. (4.24), you see that the situation is now identical to what happened in the classical case. In Eq. (4.24), we have the angular momentum term L^2/ρ^2 intact; the term $(1/2)m\Omega^2$ is a constant and plays the role of energy eigenvalue while the other term $(-E/\rho)$ is the Coulomb potential in the new radial coordinate! Everything will be fine provided the term $\hat{\mathbf{p}}_s^2/\rho$ reduces to the standard Laplacian in $D = 3$

in the ρ coordinate. If we put $d \equiv (D - 1)$, the term \hat{p}_s^2/ρ expands out to

$$\begin{aligned} \frac{1}{\rho} \frac{1}{s^d} \frac{\partial}{\partial s} \left(s^d \frac{\partial}{\partial s} \right) &= \frac{1}{\rho} \frac{1}{\rho^{d/2}} 2\sqrt{\rho} \frac{\partial}{\partial \rho} \left(\rho^{d/2} 2\sqrt{\rho} \frac{\partial}{\partial \rho} \right) \\ &= 4 \frac{1}{\rho^{(d+1)/2}} \frac{\partial}{\partial \rho} \left(\rho^{(d+1)/2} \frac{\partial}{\partial \rho} \right). \end{aligned} \quad (4.25)$$

In order for this operator to reduce to the standard Laplacian in $D = 3$, viz., $\rho^{-2} \partial_\rho (\rho^2 \partial_\rho)$ we need the condition $(1/2)(d + 1) = 2$. This gives $d = 3$ or $D = 4$.

Why $D = 4$ is special to this problem

Thus, we can map the problem of quantum isotropic oscillator in $D = 4$ to the Coulomb problem in $D = 3$. The mapping also tells you that the bound state energy of the Coulomb system is given by $E_{\text{coul}} = -(1/2)m\Omega^2$, while the parameter in the Coulomb potential $\mathcal{V}(\rho) = -Zq^2/\rho$ is given by $Zq^2 = E_{\text{osc}}$. The energy eigenvalue for the oscillator E_{osc} is given by $E_{\text{osc}} = \hbar\Omega f$ where f gives the quantization condition for the oscillator energy levels. (For a $D = 1$ oscillator, this is just $n + (1/2)$ but for $D = 4$ it is more complicated. We will comment on it later on.) Combining these two results, we find that

$$\begin{aligned} E_{\text{coul}} &= -\frac{1}{2}m\Omega^2 = -\frac{1}{2}m \left(\frac{E_{\text{osc}}}{\hbar f} \right)^2 \\ &= -\frac{1}{2}m \left(\frac{Zq^2}{\hbar f} \right)^2 = -\frac{mZ^2q^4}{2\hbar^2 f^2}. \end{aligned} \quad (4.26)$$

This allows us to find the energy eigenstates of Coulomb/Kepler problem in $D = 3$ from the energy eigenstates of the isotropic oscillator in $D = 4$. To fix f we need to deal with the angular part of the Hamiltonian with some more care (which we will discuss below) and this leads to the result that $f = f_{n_1 n_2 \ell} = (n_1 + n_2 + |\ell| + 1)$ where n_1, n_2 range over $0, 1, 2, \dots$ and $\ell = 0, \pm 1, \pm 2, \dots$. This clearly reproduces the standard hydrogen spectra.

The only problem physicists can solve is harmonic oscillator!

After all this warm up, let me show you how to model this problem rigorously [14, 15]. Since we know that the isotropic harmonic oscillator in $D = 4$ allows us to solve the problem, let us begin with a hypothetical 4-dimensional space with the coordinates (q_1, q_2, q_3, q_4) . We could introduce one radial and three angular coordinates, instead of the Cartesian coordinates q_i , in many different ways in this space. Our aim is to introduce three angular coordinates θ, ϕ and χ such that (θ, ϕ) can be mapped to the standard spherical polar angles in a $D = 3$ subspace. This requires a special coordinatization of the $D = 4$ space which is best done as follows. Pairing up the Cartesian coordinates as (q_1, q_2) and (q_3, q_4) , we can introduce two *complex* coordinates $z_1 = q_1 + iq_2$ and $z_2 = q_3 + iq_4$. We

will now introduce the radial coordinate s and three angles (θ, ϕ, χ) by the relations

$$\begin{aligned} z_1 &= q_1 + iq_2 \equiv s \cos \frac{\theta}{2} \exp \frac{i}{2}(\chi - \phi) \equiv u \exp \frac{i}{2}(\chi - \phi); \\ z_2 &= q_3 + iq_4 \equiv s \sin \frac{\theta}{2} \exp \frac{i}{2}(\chi + \phi) \equiv v \exp \frac{i}{2}(\chi + \phi). \end{aligned} \quad (4.27)$$

*Special coordinates
in $D = 4$*

The last equalities define the variables u, v which turns out to be convenient in calculations. There is a natural mapping from the complex numbers (z_1, z_2) to the standard 3-dimensional Cartesian coordinates $\mathbf{x} = (x, y, z)$. Treating z_1 and z_2 as the components of a column vector, this relation is given by

$$\mathbf{x} = z^\dagger \boldsymbol{\sigma} z, \quad (4.28)$$

where $\boldsymbol{\sigma}$ are the standard Pauli matrices. If you use the explicit form of the Pauli matrices, you will find that these relations reduce to

$$\begin{aligned} x &= z_1^* z_2 + z_1 z_2^* = 2uv \cos \phi = s^2 \sin \theta \cos \phi; \\ y &= i(z_1 z_2^* - z_1^* z_2) = iuv(e^{-i\phi} - e^{i\phi}) = s^2 \sin \theta \sin \phi; \\ z &= s^2 \cos \theta. \end{aligned} \quad (4.29)$$

This tells you that if we set $\rho = s^2$ (which we know works very well), then (ρ, θ, ϕ) is the standard spherical polar coordinates in $D = 3$.

*The oscillator
Hamiltonian in
these coordinates*

Our next job is to write down the correct Hamiltonian for the isotropic oscillator in $D = 4$ using the coordinates (s, θ, ϕ, χ) . To begin with, the metric in the 4-dimensional space, in terms of our preferred coordinates, can be easily calculated to be

$$\begin{aligned} dl^2 &= |dz_1|^2 + |dz_2|^2 = du^2 + dv^2 + \frac{1}{4}s^2(d\phi^2 + d\chi^2) + \frac{s^2}{2} \cos \theta d\chi d\phi \\ &= ds^2 + \frac{s^2}{4}(d\theta^2 + d\phi^2 + d\chi^2) + \frac{s^2}{2} \cos \theta d\chi d\phi \\ &= ds^2 + \frac{s^2}{4}(d\theta^2 + \sin^2 \theta d\phi^2) + \frac{s^2}{4}(d\chi + \cos \theta d\phi)^2. \end{aligned} \quad (4.30)$$

Therefore the kinetic energy of the particle in the 4-dimensional space is given by

$$T = \frac{1}{2}m\dot{\ell}^2 = \frac{1}{2}m \left[s^2 + \frac{s^2}{4}(\dot{\theta}^2 + \sin^2 \theta \dot{\phi}^2) + \frac{s^2}{4}(\dot{\chi} + \cos \theta \dot{\phi})^2 \right]. \quad (4.31)$$

Computing the momenta conjugate to the coordinates s, θ, ϕ, χ , we can obtain the Hamiltonian for the free particle to be

$$H_{\text{free}} = \frac{p_s^2}{2m} + \frac{2p_\theta^2}{ms^2} + \frac{2}{ms^2} \frac{(p_\phi - \cos \theta p_\chi)^2}{\sin^2 \theta} + \frac{2p_\chi^2}{ms^2}. \quad (4.32)$$

We are ultimately interested in reducing the problem to one in $D = 3$ with the momenta (p_s, p_θ, p_ϕ) . With this motivation, we re-write the above expression, after a little bit of algebra, in the form:

$$H_{\text{free}} = \frac{p_s^2}{2m} + \frac{2}{ms^2} \left(p_\theta^2 + \frac{p_\phi^2}{\sin^2 \theta} \right) + \frac{2}{ms^2 \sin^2 \theta} p_\chi (p_\chi - 2 \cos \theta p_\phi), \quad (4.33)$$

where the first two terms have the standard form familiar to us in $D = 3$. The Hamiltonian for the 4-D oscillator is obtained by adding to this the potential energy term $(1/2)m\Omega^2 s^2$. The explicit form of the operator version of this Hamiltonian will, therefore, be given by:

$$H_{\text{osc}} = \frac{1}{2m} \left[\hat{p}_s^2 + \frac{4}{s^2} \hat{L}_{\text{std}}^2 + \frac{4}{s^2 \sin^2 \theta} \left(\frac{\partial}{\partial \chi} - 2 \cos \theta \frac{\partial}{\partial \phi} \right) \frac{\partial}{\partial \chi} \right] + \frac{1}{2} m \Omega^2 s^2, \quad (4.34)$$

where \hat{L}_{std}^2 is the standard angular Laplacian on the 2-sphere and \hat{p}_s^2 is the radial part of the Laplacian.

The solution to $\hat{H}_{\text{osc}} \Psi = E_{\text{osc}} \Psi$ will now lead to the eigenfunction $\Psi(s, \theta, \phi, \chi)$ which depends on all the three angles. But we know from Eq. (4.29) that the $D = 3$ coordinates do not involve the angle χ . Therefore we shall look at the subspace of the solutions to the equation $\hat{H}_{\text{osc}} \Psi = E_{\text{osc}} \Psi$ in which Ψ is independent of χ and satisfies the constraint $\partial \Psi / \partial \chi = 0$. We see that this reduces the Hamiltonian in Eq. (4.34) to the one appearing in Eq. (4.23), except for a rescaling of the angular momentum operator by factor 4, which is of no consequence for our purpose. Rest of the analysis can now proceed exactly as we did before in the case of Eq. (4.23).

The trouble: an extra angle!

The condition $\partial \Psi / \partial \chi = 0$ translates into the requirement that the rotations in the relevant planes of the 4-dimensional space do not change the wavefunction. The angular momentum operator in the $q_1 - q_2$ plane is given by $(q_1 \partial_2 - q_2 \partial_1)$ while the angular momentum operator in $q_3 - q_4$ plane is given by $(q_3 \partial_4 - q_4 \partial_3)$. If we arrange the eigenvalues of these two operators to be equal in magnitude but opposite in signs, then we can ensure that the wavefunctions are indeed independent of the unwanted angle. So we impose the following extra condition on the 4-dimensional wavefunction:

$$(q_1 \partial_2 - q_2 \partial_1) \Psi = -(q_3 \partial_4 - q_4 \partial_3) \Psi. \quad (4.35)$$

This condition also allows us to separate the wavefunction in the 4-dimensional space to the product of two 2-dimensional oscillators in the form $\Psi = \Psi_A(q_1, q_2) \Psi_B(q_3, q_4)$ with

$$\left[\frac{\hbar^2}{2m} (\partial_1^2 + \partial_2^2) + \lambda_A - \frac{1}{2} m \Omega^2 (q_1^2 + q_2^2) \right] \Psi_A = 0, \quad (4.36)$$

and a similar equation for Ψ_B with eigenvalue λ_B . The solutions to the 2-dimensional isotropic oscillator are well known. If we take the eigenvalue of the angular momentum to be ℓ_1 then the energy eigenvalues are given by

$$\varepsilon_A(n_1, \ell_1) = \hbar\Omega(2n_1 + |\ell_1| + 1) = \frac{\hbar^2 \lambda_A(n_1, \ell_1)}{2m}, \quad (4.37)$$

and similarly for the second oscillator. But the condition in Eq. (4.35) requires us to choose $\ell_1 = -\ell_2 = \ell$, (say), so that the final solution can be written in the form

$$\Psi_{n_1, n_2, \ell} = \Psi_{A n_1, \ell}(q_1, q_2) \Psi_{B n_2, -\ell}(q_3, q_4), \quad (4.38)$$

with $\lambda = \lambda_A + \lambda_B$ being:

$$\lambda(n_1, n_2, \ell) = 4\Omega(n_1 + n_2 + |\ell| + 1) \equiv 4\Omega N. \quad (4.39)$$

Hydrogen atom in 3D = Oscillator in 4D!

This leads to the result in Eq. (4.26) with $f = f(n_1, n_2, \ell) = (n_1 + n_2 + |\ell| + 1)$ where n_1, n_2 range over $0, 1, 2, \dots$ and $\ell = 0, \pm 1, \pm 2, \dots$

Connecting up with the Runge-Lenz approach

In this approach, which maps the 3-D hydrogen atom to a 4-D isotropic oscillator, it is obvious that our system has rotational invariance in 4-dimensional space. The physical solutions, however, are restricted to those satisfying the constraint Eq. (4.35) so that the third angle does not come into the picture. This constraint is closely related to the constraint $I^2 = K^2$ which we had in the operator approach; in fact, I_i and K_i can be thought as the angular momentum operators on the relevant planes. Finally, straightforward computation will show that the wavefunctions in Eq. (4.38) are also eigenfunctions of the z -component of the Runge-Lenz vector \mathbf{M} and satisfies the relation

$$M_z \Psi_{n_1, n_2, \ell} = \left[\frac{z q^2 (n_2 - n_1)}{N} \right] \Psi_{n_1, n_2, \ell}. \quad (4.40)$$

So we have simultaneously diagonalized all the relevant conserved quantities in the approach.

Coulomb scattering is strange, too!

So far, we have been concerned with the bound state problem in the Coulomb potential. The $(1/r)$ potential introduces some conundrums in the case of scattering as well. We will conclude this chapter with a brief description of some of these issues.

Let us start by recalling the usual formalism of quantum mechanical scattering theory. The time independent Schrödinger equation in a potential $V(\mathbf{r})$ can be expressed in the form

$$(\nabla^2 + k_0^2) \psi(\mathbf{r}) = U(\mathbf{r}) \psi(\mathbf{r}) \equiv f(\mathbf{r}), \quad (4.41)$$

where

$$\frac{2m}{\hbar^2}E = k_0^2; \quad \frac{2m}{\hbar^2}V \equiv U. \quad (4.42)$$

The formal solution to Eq. (4.41) will be

$$\psi(\mathbf{r}) = \psi_0(\mathbf{r}) + (\nabla^2 + k_0^2)^{-1} f(\mathbf{r}) = \psi_0(\mathbf{r}) + \int d^3\mathbf{r}' f(\mathbf{r}') G(\mathbf{r} - \mathbf{r}'), \quad (4.43)$$

where we interpret $\psi_0(\mathbf{r})$ as an incident wave propagating towards the scattering potential and the rest (which vanishes when $V = 0$) as the scattered wave. The second equality defines the Green function for the problem which satisfies the equation $(\nabla^2 + k_0^2)G(\mathbf{r}) = \delta_D(\mathbf{r})$. Textbooks contain several formal procedures to solve this equation for the Green function but it can be done by inspection and a bit of English! We first note that everywhere except the origin, the right hand side vanishes and we have $(\nabla^2 + k_0^2)G(\mathbf{r}) = 0$; since we want an outgoing wave as the solution for a point source, we *must* have $G(r) \propto e^{ik_0 r}/r$. All we need to do is to fix the proportionality constant. Near the origin you can ignore the k_0^2 term and equation reduces to $\nabla^2 G = \delta_D(\mathbf{r})$. This is just the Poisson equation for a point particle at the origin and we know that its solution is $G = -(1/4\pi)(1/r)$ which should be the behaviour of the Green function near origin. This fixes the proportionality constant as $-(1/4\pi)$ and we get the Green function to be

Simple logic gets you the Green's function!

$$G(\mathbf{r}) = -\frac{1}{4\pi} \frac{e^{ik_0 r}}{r}. \quad (4.44)$$

If we substitute Eq. (4.44) into Eq. (4.43) we will get an integral equation for ψ because the f on the right hand side depends on ψ . One way to solve this equation is to work perturbatively, order-by-order in the potential V . To the lowest order, we plug in $\psi_0(\mathbf{r})$ — which we can take to be an incident plane wave $\exp(i\mathbf{k}_0 \cdot \mathbf{r})$ representing an incident particle with momentum $\hbar\mathbf{k}_0$. Doing this and assuming that we can approximate $|\mathbf{r} - \mathbf{r}'|^{-1} \approx (1/r)$, we can easily show that the first order solution ψ_1 is given by

$$\psi_1(\mathbf{r}) = -\frac{1}{4\pi} \frac{e^{ik_0 r}}{r} \tilde{U}(\mathbf{q}) = -\frac{1}{2\pi} \frac{m}{\hbar^2} \tilde{V}(\mathbf{q}) \frac{e^{ik_0 r}}{r}. \quad (4.45)$$

Here, $\tilde{U}(\mathbf{q})$ and $\tilde{V}(\mathbf{q})$ are the three dimensional Fourier transforms of $U(\mathbf{r})$ and $V(\mathbf{r})$ evaluated on the momentum transfer $\mathbf{q} \equiv \mathbf{k}_0 - k_0 \hat{\mathbf{r}} \equiv \mathbf{k}_i - \mathbf{k}_f$ with $\hat{\mathbf{r}}$ being the unit vector in the direction of \mathbf{r} . In the second equality, we have indicated the initial and final momentum vectors of the scattered particles as $\hbar\mathbf{k}_i$ and $\hbar\mathbf{k}_f$ respectively. Thus the Fourier transform of the scattering potential determines the lowest order correction due to scattering.

The Born approximation

In the case of spherically symmetric potential, the coefficient of $e^{ik_0 r}/r$ can depend only on the scattering angle θ between \mathbf{k}_i and \mathbf{k}_f . If we compute the current $\mathbf{j} \equiv (\hbar/m)\text{Im } \psi^* \nabla \psi$ for the outgoing wave, we find its magnitude to be

$$j_{\text{out}} = \frac{|f(\theta)|^2}{r^2} v_0 = \frac{|f(\theta)|^2}{r^2} j_{\text{in}} , \quad (4.46)$$

since $j_{\text{in}} \equiv v_0 \equiv \hbar k_0/m$. The number of particles scattered into a solid angle $d\Omega$ is given by

$$dN = j_{\text{out}} r^2 d\Omega = v_0 |f(\theta)|^2 d\Omega \equiv j_{\text{in}} \frac{d\sigma}{d\Omega} d\Omega , \quad (4.47)$$

*The scattering
cross-section*

where the last equation defines the differential scattering cross section ($d\sigma/d\Omega$). We thus get our final result for the scattering cross section in the lowest order approximation to be

$$\frac{d\sigma}{d\Omega} = |f(\theta)|^2 = \frac{1}{4\pi^2} \frac{m^2}{\hbar^4} |\tilde{V}(\theta)|^2 . \quad (4.48)$$

Let us go ahead and apply it to the Coulomb potential (which happens to be an illegal procedure on which we shall comment upon later). Using the fact that the Fourier transform of $V(\mathbf{r}) = Ze^2/r$ is $\tilde{V}(\mathbf{k}) = 4\pi Ze^2/k^2$ and the result

$$(\mathbf{k}_i - \mathbf{k}_f)^2 = 2k_0^2(1 - \cos \theta) = 4k_0^2 \sin^2 \left(\frac{\theta}{2} \right) , \quad (4.49)$$

we find that the differential cross section for the scattering in the Coulomb field is given by

$$\left(\frac{d\sigma}{d\Omega} \right) = \frac{Z^2 e^4}{(4E)^2} \text{cosec}^4 \left(\frac{\theta}{2} \right) , \quad (4.50)$$

*Familiar, but
deceptive*

with a characteristic $\text{cosec}^4(\theta/2)$ dependence. This is a very standard result (called Rutherford scattering cross section) and is, of course, described in every text book.

*Puzzle 1:
Where is \hbar ?*

But, there are a couple of things which are quite strange about this result and deserve attention. First, you notice that the result is independent of \hbar . That is a bit strange since we are supposed to be doing quantum mechanics! The cross section we have found is *exactly* what you get doing everything purely classically with no wavefunctions, no Schrödinger equations and *a la* Rutherford. It is pretty nice for Rutherford, who got a quantum result by classical analysis, but it *is* strange.

Second, the scattering problem in Coulomb potential corresponds to solving the Schrödinger equation with $E > 0$. Although rather messy,

these solutions are known and have the asymptotic form given by

$$\begin{aligned} \psi(r) \sim & \left[1 - \frac{\gamma^2}{ik(r-z)} \right] \exp[ikz + i\gamma \log k(r-z)] \\ & - \frac{\Gamma(1+i\gamma)}{\Gamma(1-i\gamma)} \left(\gamma \operatorname{cosec}^2 \frac{\theta}{2} \right) \frac{1}{kr} \exp[ikr - i\gamma \log k(r-z)] , \end{aligned} \quad (4.51)$$

where $\gamma = Ze^2/4\pi\hbar v$, $k = mv/\hbar$ and θ is the scattering angle. The first thing we notice is that the asymptotic forms of the wave are *not* of the form e^{ikr}/r . This distortion of the phase is due to the long range nature of the Coulomb field which means that everything we did above is illegal for Coulomb scattering! Next we see that one can still read off an $f(\theta)$ from the second term in Eq. (4.51). If we compute $|f(\theta)|^2$ we find that we again get the Rutherford scattering cross section. This is quite incredible because the calculation that led to Eq. (4.50) was supposed to be valid *only to the first order perturbation theory*. In writing Eq. (4.45) we did introduce an approximation, usually called the Born approximation. How come Born approximation leads to the *exact* result for the scattering cross section? What do all the higher order (“unBorn”) terms contribute?

Puzzle 2:
How can Born approximation give the exact result?

The answer to the first puzzle is relatively simple but the second one is more involved. We can understand why there is no \hbar in the final result by the following scaling argument. If $V(r) \sim r^n$ then $\tilde{V}(k) \sim k^{-(3+n)}$. Therefore,

$$|f|^2 \sim \frac{1}{\hbar^4} k^{-2(3+n)} \sim \frac{1}{\hbar^4} \frac{\hbar^{2(n+3)}}{(\hbar k)^{2(n+3)}} \sim \frac{1}{E^{(n+3)}} \hbar^{2n+2} , \quad (4.52)$$

leading to

$$\left(\frac{d\sigma}{d\Omega} \right) \propto \frac{\hbar^{2(n+1)}}{E^{(n+3)}} . \quad (4.53)$$

Once again we see the special status enjoyed by the Coulomb potential with $n = -1$. *This is the only power law potential for which the scattering cross section is independent of \hbar just because of dimensional reasons.*

Once again, inverse-square is special!

To understand the second issue, we actually need to compute the higher order terms beyond Born approximation and see what they do. This has been done in the literature (see, for example, Ref. [16]). To do things in a well defined manner, one can calculate the scattering cross section order-by-order for a screened Coulomb potential of the form $e^{-\lambda r}/r$ and then take the limit of $\lambda \rightarrow 0$. Such a calculation shows that all the higher order terms only change the *phase* of the outgoing scattered wave leaving $|f(\theta)|^2$ invariant. Unfortunately, no one knows a simple reason as to why this happens — which makes it an interesting question for further exploration.

A strange, but calculable, result

Potential surprises in Newtonian Gravity

5

Consider a planet which has a weird shape, resembling, say, that of a diseased potato. Is it possible that the gravitational force exerted by this planet — which is distinctly non-spherical in shape — falls exactly as r^{-2} *everywhere* outside of it? The initial reaction of many physicists will be: “No, of course, not; you need a spherically symmetric distribution of mass to produce a $1/r^2$ force outside it”. Surprisingly, this is not true. You can construct totally weird mass distributions which exert an inverse square law force on the outside world.

Incredible, but true!

To begin with, let me assure you that there is no cheating involved here. We are not talking about the gravitational field far away from the body which falls *approximately* as $1/r^2$. The result should be *exact* and must hold everywhere outside the the body, right from its surface. You also need not worry about things like viewing a spherically symmetric distribution in a strange coordinate system etc.. We are thinking of standard Cartesian coordinates with concepts like spherical symmetry having the usual meaning.

To understand the implications of the question, we start by reviewing some basics of Newtonian gravity. The Newtonian gravitational field \mathbf{F} can be expressed as the gradient of a potential ϕ which satisfies Poisson equation. We have

$$\nabla^2 \phi = 4\pi G\rho; \quad \mathbf{F} = -\nabla \phi, \quad (5.1)$$

where $\rho(\mathbf{x})$ is the matter density which is assumed to be either positive or zero everywhere. (For our purpose, it is adequate to consider static configurations.) In these mathematical terms our problem translates to the question: Can you find a density distribution $\rho(\mathbf{x})$ which is *not* spherically symmetric (in some chosen coordinate system) and vanishes outside some compact region \mathcal{R} around the origin, such that outside \mathcal{R} the potential ϕ falls as $1/r$? Of course, any spherically symmetric $\rho(\mathbf{x})$ will produce such a potential, but *must* it be spherically symmetric?

*Question, made
precise*

*More general
question*

A little thought will convince you that there is no simple way of going about analyzing this problem. Usually, we are given some $\rho(\mathbf{x})$ and asked to find the $\phi(\mathbf{x})$. We are now interested in the inverse question, which — in a broader context — is the following: If we know the gravitational force in some region of space how unique is the density distribution producing that force? (Some of these issues are discussed in classical, geometrical style in older books on potential theory, like e.g., Refs. [17, 18]; also see [19].)

Let me give you some instances in which totally different density distributions produce the same gravitational field in some region. This will be a good warm up for the original question we want to attack.

One example, well-explored in standard text books, is the field produced by an infinite, plane sheet of matter of surface mass density σ . You might not have learnt it in the context of gravity, but I am sure you have encountered it in some electrostatics course. You will remember that such infinite planes with constant surface density produce a gravitational force $\mathbf{F} = -2\pi G\sigma\hat{\mathbf{n}}$ which is constant everywhere and directed towards the sheet. (Here, $\hat{\mathbf{n}}$ is the unit vector in the direction perpendicular to the sheet.) We now ask: Is it possible to come up with a density distribution which is not plane-symmetric but will produce *constant* gravitational field in some compact region of space \mathcal{S} ? The answer is “yes”; and some of you must have even worked it out without quite realizing its importance!

*Virtues of
superposition*

The configuration is shown in Fig 5.1. Consider a sphere, of radius R and constant density ρ , centered at the origin of the coordinate system. Inside it we carve out another spherical region of radius L centered at the point ℓ . Consider the force on a particle located within the cavity at the position $(\ell + \mathbf{r})$. The force due to a constant density sphere is $\mathbf{F} = -(4/3)\pi G\rho\mathbf{x}$ (so that $-\nabla\cdot\mathbf{F} = 4\pi G\rho$). Hence, the force we want is

$$\mathbf{F} = \mathbf{F}_{sph} - \mathbf{F}_{hole} = -\frac{4}{3}\pi G\rho[\ell + \mathbf{r}] + \frac{4}{3}\pi G\rho\mathbf{r} = -\frac{4\pi}{3}G\rho\ell, \quad (5.2)$$

where \mathbf{F}_{hole} is the force the matter in cavity would have exerted if it were not empty. This \mathbf{F} is clearly a constant inside the hole! Thus a spherical hole (located off-center) inside a sphere is a region with constant gravitational force!

Suppose you measure the gravitational field in some finite region \mathcal{S} and find it to be strictly constant. Can you say anything about the mass distribution which is producing this force? Of course not. It could have been produced by an infinite plane sheet or a hole-in-a-sphere; these are just two of infinitely many possibilities. Most of these mass distributions, which produce a constant gravitational field, will not have any specific symmetry.

*Zero-gravity is easy,
when you have
constant gravity*

We can twist around the hole-in-the-sphere example to lead to another interesting conclusion. You must have learnt, while studying Newtonian gravity, that a spherical shell of matter exerts no gravitational force on

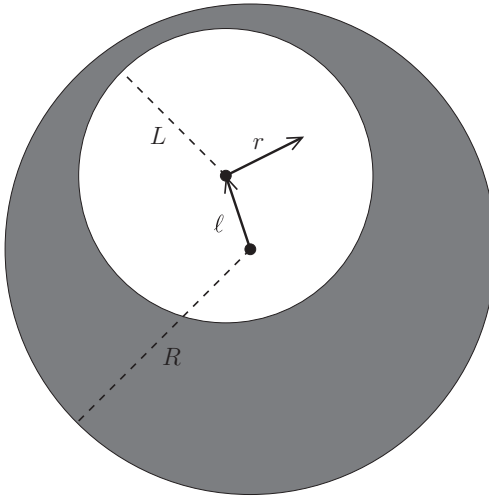


Fig. 5.1: An example of a highly asymmetric density distribution leading to a constant gravitational force in a compact region of space. We scoop out all matter from a spherical region of radius L located inside a constant density distribution which originally made a sphere of radius R . It is easy to show that everywhere inside the spherical hole the gravitational force is a constant and is in the direction along the vector joining the centers of the two spheres.

a particle inside it. (This is just a special case of Eq. (5.2) above; when $\ell = 0$, the force vanishes.) Is it possible to come up with a completely asymmetric distribution of matter which exerts *zero* gravitational force in some region?

It turns out that the answer is again ‘yes’ and all you need to do is the following: Suppose you make two hole-in-the-sphere distributions with different values for the parameters — one with density ρ_1 , radius R_1 , hole radius L_1 and the center of the hole located at ℓ_1 with respect to the center of the sphere; the second one has density ρ_2 , radius R_2 etc. We superpose the spheres such that: (i) ℓ_1 and ℓ_2 are in the opposite directions; (ii) $\rho_1 \ell_1 = \rho_2 \ell_2$; and (iii) part of the spherical cavities overlap. The resulting density distribution is clearly not spherically symmetric. But in the region of the cavity which is common to the holes of both spheres, the gravitational force is strictly zero. This is because each sphere produces an equal and opposite force in the cavity when $\rho_1 \ell_1 = \rho_2 \ell_2$.

The moral of the story is worth remembering. Just knowing the symmetries of the gravitational force in *some* region alone does not allow you to conclude anything about the *symmetries* of the mass distribution. This itself comes as a surprise for many physicists since we are so accustomed to assuming the same symmetries for the field and its source. It is quite possible to have completely asymmetric density distributions producing

Sources and fields need not share the same symmetry in finite regions

highly regular gravitational fields. We are now ready to tackle the question we originally started with: Are there density distributions which are *not* spherically symmetric but produce an inverse square force?

Let us begin by considering this problem in the case of electrostatics. Is it possible to have a charge distribution which is not spherically symmetric but produces an inverse square electric field? Incredibly enough, you already know such a distribution from your regular electrostatics course! Remember the problem of a point charge and a conducting sphere which is solved by the method of images? We start with a conducting sphere of radius a and a point charge Q located outside the sphere at a distance L from the center of the sphere. The charge Q induces a surface charge distribution on the conducting sphere and the net electric field at any point is the sum of the electric fields due to the surface charge distribution σ and the point charge Q . This problem is solved by showing that it is equivalent to that of two point charges: the real charge Q and an “image” charge $q = -(a/L)Q$ placed at a distance $\ell = (a^2/L)$ inside the sphere in the line joining the center of the sphere to the charge Q . The fields outside this sphere, produced by the point charges Q and q , are identical to those due to the point charge Q and the charge distribution σ . It follows that this charge distribution σ produces a field which is equivalent to that of a point charge q ! The explicit form of the charge distribution is given by

$$\sigma(r, \theta) = -\frac{Q}{4\pi r} \frac{(L^2 - r^2)}{(L^2 + r^2 - 2Lr \cos \theta)^{3/2}}; \quad (\text{at } r = a). \quad (5.3)$$

Of course, this distribution σ is far from spherically symmetric since the induced charge on the side nearer to Q will be distributed differently compared to the induced charge on the farther side. We have thus come up with a charge distribution which is not spherically symmetric but produces a strict inverse square law force outside a finite region.

The main difference between electrostatics and gravity is that, in gravity, the mass density has to be positive definite — while, in electrostatics, the charge density need not be positive definite. In the above example, the charge density has the same sign everywhere and hence one can simply replace it by mass density to get a solution appropriate for the gravitational case.

If you are still shaking your head in disbelief, let me assure you that everything is quite above board. It is quite possible to have such distributions, and — in fact — there are infinitely many such configurations. Those of you who are mathematically inclined might like the following construction of some such distributions using a property of Poisson equation known as “inversion”. Inversion is a mathematical operation under which you associate to any point \mathbf{x} another point $\mathbf{x}_{inv} \equiv (a^2/x^2)\mathbf{x}$ where a is the radius of the “inverting sphere”. From this definition it immediately follows that points inside a sphere of radius a are mapped to points out-

*Actually you
already know this*

*Non-spherical
charge density
leading to $1/r^2$
electric field!*

*Poisson equation
under inversion
— a result worth
knowing*

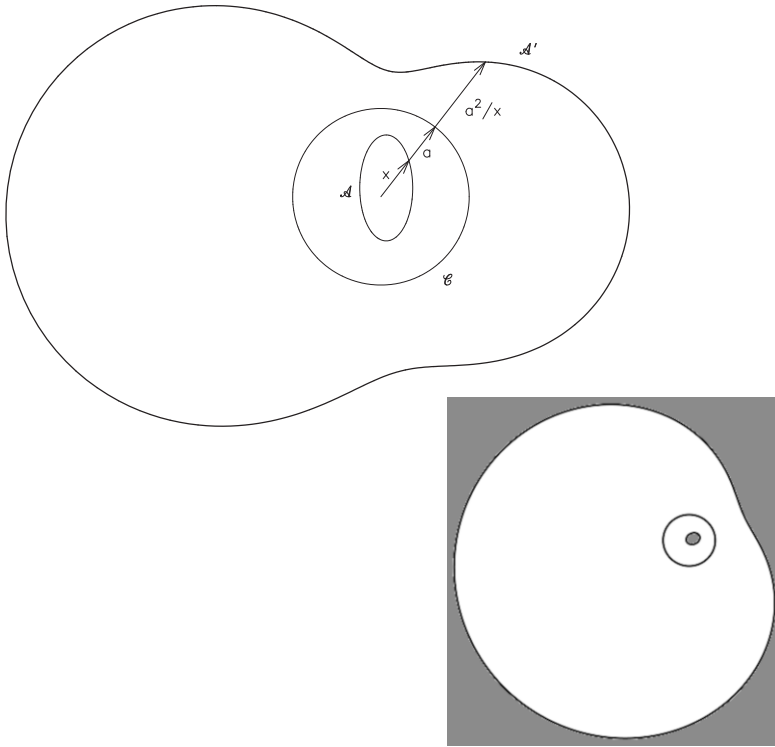


Fig. 5.2: *Top left:* Schematic picture showing the effect of inversion in which a given point \mathbf{x} is mapped to another point $\mathbf{x}_{inv} \equiv (a^2/x^2)\mathbf{x}$. When the points in a surface of a compact region \mathcal{A} are inverted using the inverting sphere \mathcal{C} , we obtain the surface \mathcal{A}' . In this process, the region inside \mathcal{A} gets mapped to region outside \mathcal{A}' and vice versa. *Bottom right:* Actual inversion of a black, shaded, oval shaped region by a sphere. The inverted curve is shown with shaded region being mapped to the outside.

side and vice-versa. There is an interesting connection between inversion and the solutions to the Poisson equation. Suppose $\phi[\mathbf{x}; \rho(\mathbf{x})]$ is the gravitational potential at a point \mathbf{x} due to a density distribution $\rho(\mathbf{x})$. Consider now a new density distribution $\rho'(\mathbf{x}) = (a/x)^5 \rho(\mathbf{x}_{inv})$ obtained by taking the original density at the inverted point $\mathbf{x}_{inv} \equiv (a^2/x^2)\mathbf{x}$ and multiplying by $(a/x)^5$. We can show that the gravitational potential due to $\rho'(\mathbf{x})$ is given by $\phi'(\mathbf{x}) = (a/x)\phi(\mathbf{x}_{inv})$. That is, the new gravitational potential at any given point is the old gravitational potential at the inverted point \mathbf{x}_{inv} multiplied by (a/x) . (You will find the proof in the Appendix at the end of this chapter.)

This result can be used to produce strange looking compact mass distributions with strictly inverse square law force. We begin with the result, obtained earlier, that one can have very asymmetric density distributions which can produce *zero* gravitational force inside an empty compact re-

gion of space \mathcal{A} . In figure 5.2, we assume that there are sources outside of \mathcal{A} (which are not shown) that produce a constant gravitational potential inside the region of space \mathcal{A} . The exact shape of this region is immaterial for our discussion.

Let \mathcal{C} be an imaginary spherical surface of radius a with center somewhere inside \mathcal{A} . We now invert the surface of the region \mathcal{A} using the inverting sphere \mathcal{C} and obtain the surface \mathcal{A}' . In this process, the compact region inside \mathcal{A} gets mapped to an infinite, non-compact region outside \mathcal{A}' . Since the region inside \mathcal{A} was originally empty, the region outside \mathcal{A}' will be empty in the inverted configuration; all the sources which were originally outside \mathcal{A} are now mapped to the region inside \mathcal{A}' . Consider now the gravitational potential outside \mathcal{A}' due to this source ρ' which is now inside \mathcal{A}' . This potential is obtained by taking the potential due to the inverted point inside \mathcal{A} and multiplying it by (a/x) . But since the potential everywhere inside \mathcal{A} is a constant it follows that the potential outside \mathcal{A}' falls as $|\mathbf{x}|^{-1}$. We now have a region \mathcal{A}' outside which the gravitational force is strictly inverse square and the density distribution producing this force is far from spherical!

Turning it inside out

Historically, this problem seems to have been first raised (and answered) by Lord Kelvin. Newton, on the other hand, never worried about this question. This is somewhat surprising since Newton had worried a lot about the related problem, viz., whether a spherically symmetric mass distribution will produce a force as though all its mass is concentrated at the origin.

Appendix: In the text we used a result connecting the Newtonian gravitational potential to its source when we perform an inversion. The proof of this result is outlined here. A clever way to prove this result is to consider the effect of conformal transformations of the Laplace equations which is outlined in Ref. [20], page 151. A more elementary procedure, though algebraically involved, is as follows.

Starting from the Poisson equation $\nabla^2\phi = 4\pi G\rho$ relating the gravitational potential ϕ to matter density ρ , we write down the solution:

$$-\phi(\mathbf{x}) = G \int \frac{\rho(\mathbf{r})}{|\mathbf{x} - \mathbf{r}|} d^3\mathbf{r}. \quad (5.4)$$

We have the vector identity for any \mathbf{r}, ℓ which reads:

$$\frac{1}{|\mathbf{r} - (a^2/\ell^2)\ell|} = \frac{|\ell|}{|\mathbf{r}|} \frac{1}{|(a^2/r^2)\mathbf{r} - \ell|}. \quad (5.5)$$

Identifying \mathbf{x} with $(a^2/\ell^2)\ell$ we can write

$$-\phi\left(\frac{a^2}{\ell^2}\ell\right) = G \int \frac{\rho(\mathbf{r})d^3\mathbf{r}}{|\mathbf{r} - (a^2/\ell^2)\ell|} = G|\ell| \int \frac{\rho(\mathbf{r})|\mathbf{r}|}{|(a^2/r^2)\mathbf{r} - \ell|} d^3\mathbf{r}. \quad (5.6)$$

We transform the (dummy) integration variable \mathbf{r} to \mathbf{R} , with $\mathbf{r} = (a^2/R^2)\mathbf{R}$; $d^3\mathbf{r} = (a^6/|R|^6)d^3\mathbf{R}$, getting:

$$\begin{aligned} -\phi\left(\frac{a^2}{\ell^2}\boldsymbol{\ell}\right) &= G|\boldsymbol{\ell}| \int \frac{a^6}{R^6} \frac{R}{a^2} \frac{\rho((a^2/R^2)\mathbf{R})}{|\mathbf{R}-\boldsymbol{\ell}|} d^3\mathbf{R} \\ &\equiv \left(\frac{\ell}{a}\right) G \int \frac{\eta(\mathbf{R})d^3\mathbf{R}}{|\mathbf{R}-\boldsymbol{\ell}|} \equiv -\frac{|\boldsymbol{\ell}|}{a} u(\boldsymbol{\ell}) , \end{aligned} \quad (5.7)$$

where $u(\boldsymbol{\ell})$ is the potential due to $\eta(\mathbf{x})$. That is, $\nabla^2 u = 4\pi G\eta$. This gives the relation between potential-density pairs of the form:

$$\phi\{\mathbf{x}; \rho(\mathbf{x})\} = \frac{a}{|\mathbf{x}|} \phi\left\{\frac{a^2}{x^2}\mathbf{x}; \frac{a^5}{x^5}\rho\left(\frac{a^2}{x^2}\mathbf{x}\right)\right\}. \quad (5.8)$$

Potential at \mathbf{x} due to a distribution $\rho(\mathbf{x})$ is the same as $(a/|\mathbf{x}|)$ times the potential at $(a^2/x^2)\mathbf{x}$ due to a distribution $(a^5/x^5)\rho((a^2/x^2)\mathbf{x})$. This is the result we used in the text.

The idealized problem of a planet orbiting around the Sun has an exact solution which — as we saw in Chapter 3 — is fairly easy to obtain. But in real life, the orbital motion of planets is a lot more complicated because each planet is influenced by the gravitational force of all other bodies in the solar system. In fact, if we add just one more gravitating body — thereby reaching the three-body problem, in which three point particles are moving under the gravitational influence of one another — the problem becomes analytically intractable.

When an exact problem cannot be solved, physicists attempt to solve a simpler version of the problem, which will at least capture some features of the original one. One such case corresponds to what is known as the *restricted three-body problem* which could be described as follows. Consider two particles of masses m_1 and m_2 which orbit around their common center of mass exactly as in the case of the standard Kepler problem. We now consider a third particle of mass m_3 , with $m_3 \ll m_1$ and $m_3 \ll m_2$, in the gravitational field of the two particles m_1 and m_2 . Since it is far less massive than the other two particles, we will assume that it behaves like a test particle and does not affect the original motion of m_1 and m_2 . You can see that this is equivalent to studying the motion of m_3 in a time dependent external gravitational potential produced by the masses m_1 and m_2 . Given the fact that we have lost both the time translation invariance and axial symmetry, any hope for simple analytic solutions is misplaced. But there is a special case for which a truly beautiful solution can be obtained.

*Tractable
version the of
3-body problem*

This corresponds to a situation in which all the three particles maintain their relative positions with respect to one another but rotate rigidly in space with an angular velocity ω ! In fact, the three particles are located at the vertices of an equilateral triangle irrespective of the ratio of the masses m_1/m_2 . If you think about it, you will find that this solution, first found by Lagrange, is not only elegant but also somewhat counter-intuitive. How are the forces, *which depend on mass ratios*, balanced without adjusting the distance ratios but always maintaining the equilateral configuration?

*A surprisingly
elegant solution!*

*Stable orbits around
potential maxima!*

What is more, the location of m_3 happens to be at the local *maximum* of the effective potential in the frame co-rotating with the system. Traditionally, the maxima of a potential have a bad press due to their tendency to induce instability. It turns out that, in this solution, stability can be maintained (for a reasonable range of parameters) because of the existence of Coriolis force — which is one of the concepts for which it is difficult to acquire an intuitive grasp. I will now derive this solution and describe its properties [21].

If the separation between the masses m_1 and m_2 is a , the Kepler solution implies that they can rotate in circular orbits around the center of mass with the angular velocity given by

$$\omega^2 = \frac{G(m_1 + m_2)}{a^3}, \quad (6.1)$$

*Rotate with the
masses*

where a is the distance between the particles. Since Lagrange has told us that a rigidly rotating solution exists with the third body, we will study the problem in the coordinate system co-rotating with the masses in which the three bodies are at rest. We will first work out the equations of motion for a particle in a rotating frame before proceeding further.

This is most easily done by starting from the Lagrangian for a particle $L(\mathbf{x}, \dot{\mathbf{x}}) = (1/2)m\dot{\mathbf{x}}^2 - V(\mathbf{x})$ and transforming it to a rotating frame, by using the transformation law $\mathbf{v}_{\text{inertial}} = \mathbf{v}_{\text{rot}} + \boldsymbol{\omega} \times \mathbf{x}$ where $\boldsymbol{\omega}$ is the angular velocity of the rotating frame. Substituting into L leads to the Lagrangian of the form

$$L = \frac{1}{2}m\mathbf{v}^2 + m\mathbf{v} \cdot (\boldsymbol{\omega} \times \mathbf{x}) + \frac{1}{2}m(\boldsymbol{\omega} \times \mathbf{x})^2 - V(\mathbf{x}); \quad \mathbf{v} \equiv \mathbf{v}_{\text{rot}}, \quad (6.2)$$

*Putting the
Lagrangian
to good use*

and the corresponding equations of motion will be:

$$m \frac{d\mathbf{v}}{dt} = -\frac{\partial V}{\partial \mathbf{x}} + 2m\mathbf{v} \times \boldsymbol{\omega} + m\boldsymbol{\omega} \times (\mathbf{x} \times \boldsymbol{\omega}). \quad (6.3)$$

We see that the transformation to a rotating frame introduces two additional force terms in the right hand side of Eq. (6.3), of which, the $2m(\mathbf{v} \times \boldsymbol{\omega})$ is called the Coriolis force and $m\boldsymbol{\omega} \times (\mathbf{x} \times \boldsymbol{\omega})$ is the more familiar centrifugal force. The Coriolis force has a form identical to the force exerted by a magnetic field $(2m/q)\boldsymbol{\omega}$ on a particle of charge q . It follows that this force cannot do any work on the particle since it is always orthogonal to the velocity. The centrifugal force, on the other hand, can be obtained as the gradient of an effective potential which is the third term on the right hand side of Eq. (6.2).

*This is done
with foresight*

We can now find the solution to the rigidly rotating system, in which all the three particles are at rest in the rotating frame in which Eq. (6.3) holds. We will choose a coordinate system in which *the test particle is at the origin* and denote by $\mathbf{r}_1, \mathbf{r}_2$ the position vectors of masses m_1 and m_2 .

The position of the center of mass of the m_1 and m_2 will be denoted by \mathbf{r} , so that:

$$(m_1 + m_2)\mathbf{r} = m_1\mathbf{r}_1 + m_2\mathbf{r}_2. \quad (6.4)$$

For the solution we are looking for, all these three vectors are independent of time in the rotating frame and the Coriolis force vanishes because $\mathbf{v} = 0$. Since the rotational motion of m_1 and m_2 is already taken care of (and they are assumed to be oblivious to m_3), we only need to satisfy the equation of motion for m_3 . This demands:

$$\frac{Gm_1}{r_1^3}\mathbf{r}_1 + \frac{Gm_2}{r_2^3}\mathbf{r}_2 = \omega^2\mathbf{r}. \quad (6.5)$$

You should now be able to see the equilateral triangle emerging. If we assume $r_1 = r_2$, and take note of Eq. (6.4), the left hand side of Eq. (6.5) can be reduced to $(G/r_1^3)(m_1 + m_2)\mathbf{r}$ which is in the direction of \mathbf{r} . If we next set $r_1 = a$, Eq. (6.5) is identically satisfied, thanks to Eq. (6.1). (The cognoscenti would appreciate the algebraically clever trick of making the location of the test particle as the origin.) This analysis shows how the mass ratios go away through the proportionality of both sides to the radius vector between the center of mass and the test particle.

The equilateral triangle

To ensure that we obtain *all* the equilibrium solutions, we can do this more formally. If we define the vector \mathbf{q} by the relation $m_1\mathbf{r}_1 - m_2\mathbf{r}_2 = (m_1 + m_2)\mathbf{q}$, a little bit of algebraic manipulation allows us to write Eq. (6.5) as:

$$\frac{G(m_1 + m_2)}{2r_1^3r_2^3}[(r_1^3 + r_2^3)\mathbf{r} + (r_2^3 - r_1^3)\mathbf{q}] = \frac{G(m_1 + m_2)}{a^3}\mathbf{r}. \quad (6.6)$$

For this equation to hold, all the vectors appearing in it must be collinear. One possibility is to have \mathbf{r} and \mathbf{q} to be in the same direction. It then follows that $\mathbf{r}_1, \mathbf{r}_2$ and \mathbf{r} are all collinear and the three particles are in a straight line. The equilibrium condition can be maintained at three locations, usually called L_1, L_2 and L_3 . To work out the exact position of equilibrium, one has to solve a fifth-order equation which will lead to three real roots. We are, however, not interested in these, though L_2 of the Sun-Earth system has lots of practical applications.

Lagrange has five points to make

If we do *not* want \mathbf{r} and \mathbf{q} to be parallel to each other, then the only way to satisfy Eq. (6.6) is to make the coefficient of \mathbf{q} vanish which requires $r_1 = r_2$. Substituting back into Eq. (6.6), we find that each should be equal to a . So we get the rigidly rotating equilateral configuration of three masses with:

$$r_1 = r_2 = a. \quad (6.7)$$

Obviously, there are two such configurations corresponding to the two equilateral triangles we can draw with the line joining m_1 and m_2 as one side. The locations of the m_3 corresponding to these two solutions are called L_4 and L_5 , giving Lagrange a total of five points.

*Nature, of course,
knows this solution*

Incidentally, there are several examples in the solar system where we find nature using Lagrange's insight. The most famous among them is the collection of thousands of asteroids called the Trojans which are located at the vertex of an equilateral triangle, the base of which is formed by the Sun and Jupiter — the two largest gravitating bodies in the solar system. (See Box 6.1.)

The existence of such real life solutions tells us that the equilateral solution must be stable in the sense that if we displace m_3 from the equilibrium position L_5 slightly, it should come back to it. (It turns out that the other three points L_1, L_2, L_3 are *not* stable, which is easy to prove.) Our next task is to study this stability, for which a different coordinate system is better [22]. We will now take the origin of the rotating coordinate system to be *at the location of the center of mass* of m_1 and m_2 with the x-axis passing through the two masses, and the motion confined to the x-y plane.

Box 6.1: The Trojans (and the Greeks)

The solar system is replete with examples of nature making use of the Lagrange points L_4 and L_5 . The classic case is that of over 850 so called Trojan asteroids which form an equilateral triangle with the Jupiter – Sun system. In addition, the Saturn – Sun system has a few, the Mars – Sun system has two and the Neptune – Sun system has about five. Due to various other perturbing effects, some of the “Trojans” are expected to escape from the bound state within the age of the solar system. So, occasionally, they pose a bit of theoretical puzzle in planetary dynamics.

*The Greeks and
the Trojans, up in
the sky*

The first Trojan asteroid of the Sun – Jupiter system was discovered by Max Wolf in 1906 and named Achilles. The asteroids discovered subsequently in Jupiter's Lagrangian points were all given names associated with the heroes in the Iliad. Just to be fair to both sides in the Trojan war, those at the L_4 point are named after the Greek heroes and those at the L_5 point are named after the heroes of Troy. Unfortunately, the first one discovered at the L_5 point was called Patroclus (a Greek) before the Greece-Troy rule was devised. Thus a Greek name appears in the Trojan side; however, as though to compensate, Hector, the Trojan hero appears in the Greek side (and is also the largest of the Trojan asteroids). Except for these, the two sides are well segregated. Right now the Greeks (4021) outnumber Trojans (2052) nearly two-to-one! (The list of minor planets can be found at the website: <http://www.minorplanetcenter.org/iau/lists/JupiterTrojans.html>)

It will also help to rescale the variables to simplify life. Measuring all the masses in terms of the total mass $(m_1 + m_2)$, we can denote the smaller mass by μ and the larger by $(1 - \mu)$. Similarly, we will measure all distances in terms of the separation a between the two primary masses and choose the unit of time such that $\omega = 1$. The position of m_3 is (x, y) while r_1 and r_2 will denote the (scalar) distances to m_3 from the masses $(1 - \mu)$ and μ respectively (Note that these are *not* the distances to m_3 from the origin.). It is now easy to see that the equations of motion, given by Eq. (6.3), reduce to the set:

Such tricks are worth learning.

$$\ddot{x} - 2\dot{y} = -\frac{\partial \Phi}{\partial x}, \quad \ddot{y} + 2\dot{x} = -\frac{\partial \Phi}{\partial y}, \quad (6.8)$$

where

$$\Phi = -\frac{1}{2}(x^2 + y^2) - \frac{(1 - \mu)}{r_1} - \frac{\mu}{r_2} \quad (6.9)$$

is the effective potential in the rotating frame. The first term in Eq. (6.9) gives rise to the centrifugal force while the other two terms are the standard gravitational potential energy. The only known integral of motion to this equation is the rather obvious one corresponding to the energy function $(1/2)v^2 + \Phi = \text{constant}$. A little thought shows that $\nabla \Phi = 0$ at L_4 and L_5 confirming the existence of a stationary solution. To study the stability, we normally would have checked whether these correspond to a maxima or minima of the potential. As we can see from Fig. 6.1, the L_4 and L_5 actually correspond to maxima of Φ , so, if that is the whole story, L_4 and L_5 should be unstable.

Stability at the maxima of the potential?!

But, of course, that is not the whole story since we need to take into account the Coriolis force term corresponding to $(2\dot{y}, -2\dot{x})$ in Eq. (6.8). To see the effect of this term clearly, we will write the Coriolis force term in Eq. (6.3), as $(C\dot{y}, -C\dot{x})$, so that the real problem corresponds to $C = 2$. But this trick allows us to study the stability for any value of C , including for $C = 0$, to see what happens if there is no Coriolis force. We now have to do a Taylor series expansion of the terms in Eq. (6.8) in the form $x(t) = x_0 + \Delta x(t)$, $y(t) = y_0 + \Delta y(t)$ where the point (x_0, y_0) corresponds to the L_5 point with $y_0 > 0$. We also need to expand Φ up to quadratic order in Δx and Δy to get the equations governing the small perturbations around the equilibrium position. This is straightforward but a bit tedious. If you work it through, you will get the equations

Another trick: switching the Coriolis force on and off!

$$\frac{d^2}{dt^2}\Delta x = \frac{3}{4}\Delta x + \left(\frac{3\sqrt{3}}{4}\right)(1 - 2\mu)\Delta y + C\frac{d}{dt}\Delta y; \quad (6.10)$$

$$\frac{d^2}{dt^2}\Delta y = \frac{9}{4}\Delta y + \left(\frac{3\sqrt{3}}{4}\right)(1 - 2\mu)\Delta x - C\frac{d}{dt}\Delta x. \quad (6.11)$$

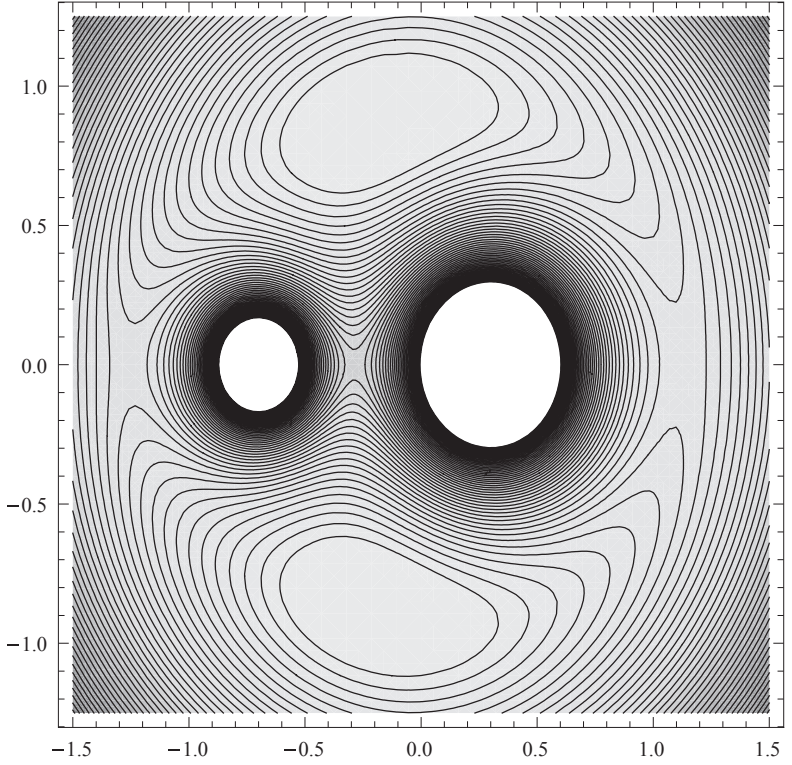


Fig. 6.1: A contour plot of the potential $\Phi(x, y)$ when $\mu = 0.3$. The L_4 and L_5 are at the potential maxima. One can also see the saddle points L_1, L_2, L_3 along the line joining the two primary masses.

To check for stability, we try solutions of the form $\Delta x = A \exp(\lambda t)$, $\Delta y = B \exp(\lambda t)$ and solve for λ . An elementary calculation gives:

$$\lambda^2 = \frac{3 - C^2 \pm [(3 - C^2)^2 - 27\mu(1 - \mu)]^{1/2}}{3}. \quad (6.12)$$

Stability requires that we should *not* have a positive real part to λ ; that is, λ^2 must be real and negative. For λ^2 to be real, the term in Eq. (6.12) containing the square root should have a positive argument. This requires

$$(C^2 - 3)^2 > 27\mu(1 - \mu). \quad (6.13)$$

Further, if both roots of λ^2 are negative, then the product of the roots must be positive and the sum must be negative. It is easily seen that this requires the condition $C > \sqrt{3}$. Hence we conclude that the motion is unstable if $C < \sqrt{3}$; in particular, in the absence of the Coriolis force ($C = 0$), the

motion is unstable because the potential at L_5 is actually a maximum. But when $C > \sqrt{3}$ — and in particular for the real case we are interested in with $C = 2$ — the motion is stable when the condition in Eq. (6.13) is satisfied. Using $C = 2$, we can reduce this condition to $\mu(1 - \mu) < (1/27)$. This leads to

$$\mu < \left(\frac{1}{2} - \sqrt{\frac{23}{108}} \right) \approx 0.0385. \quad (6.14)$$

This criterion is met by the Sun-Jupiter system with $\mu \approx 0.001$ and by the Earth-Moon system with $\mu \approx 0.012$. The stability of the Trojans is assured. In fact, L_5 and L_4 are favourites of science fiction writers and some NASA scientists for setting up space colonies. (There is even a US based society called the L_5 society, which was keen on space colonization based on L_5 !)

All fine with Jupiter and the Moon

The algebra is all fine but how does Coriolis force *actually* stabilize the motion? When the particle wanders off the maxima, it acquires a non-zero velocity and the Coriolis force induces an acceleration in the direction perpendicular to the velocity. As we noted before, this is just like motion in a magnetic field and the particle just goes around L_5 . The idea that a force which does not do work, can still help in maintaining the stability, may appear a bit strange but is completely plausible. In fact, the analogy between the Coriolis and magnetic forces tells you that one may be able to achieve similar results with magnetic fields too. This is true (and one example is the so called Penning trap).

What happens when a Trojan wanders off?

To be absolutely correct — and for the sake of experts who may be reading this — I should add a comment regarding another peculiarity which this system possesses. A more precise statement of our result on stability is that, when Eq. (6.14) is satisfied, the solutions are *not linearly unstable*. The characterization “not unstable” is qualified by saying that this is a result in linear perturbation theory. A more complex phenomenon (which is too sophisticated to be discussed here, but see Ref. [23] if you are interested) makes the system unstable for two precise values of μ which do satisfy Eq. (6.14). These values happen to be $(1/30)[15 - \sqrt{213}]$ and $(1/90)[45 - \sqrt{1833}]$. (Yes, but I did say that the phenomenon is complex!) While this is of great theoretical value, it is not of much practical relevance since one cannot fine-tune masses to any precise values.

Comment for the fussy expert

In 1697, Bernoulli announced a challenge to the mathematicians with the following words: “I, Johann Bernoulli, greet the most clever mathematicians in the world. Nothing is more attractive to intelligent people than an honest, challenging problem whose possible solution will bestow fame and remain as a lasting monument. Following the example set by Pascal, Fermat, etc., I hope to earn the gratitude of the entire scientific community by placing before the finest mathematicians of our time a problem which will test their methods and the strength of their intellect. If someone communicates to me the solution of the proposed problem, I shall publicly declare him worthy of praise”.

In those days, they did it differently

The problem he proceeded to pose was known as the brachistochrone problem (*brachistos* meaning shortest and *chronos* referring to time) which requires us to find a curve connecting two points *A* and *B* in a vertical plane such that a bead, sliding along the curve under the action of gravity, will travel from *A* to *B* in the shortest possible time.

It was known to Johann Bernoulli (and to several others, see Box 7.1 for a taste of history) that this curve is (a part of) a cycloid if we take the Earth’s gravitational field to be constant. The cycloidal path *also* has the property that time taken for a particle to roll from any point to the minima of the curve is independent of where it started from. In other words, a particle executing oscillations in a cycloidal track under the action of gravity will maintain a period which is independent of amplitude. This is quite valuable in the construction of pendulum clocks and the early clock makers knew this well. (This earned the cycloid the names isochrone and tautochrone, as though brachistochrone was not a mouthful enough!)

Cycloid: Solution to all chronic problems?

The cycloid is the curve traced by a point on the circumference of a wheel, which rolls without slipping, along a straight line. It is easy to show (see Fig. 7.1; the figures for cycloids in some published works are incorrect in the sense that the tangents at the extremities make arbitrary angles with the axis!) that the parametric equation ($x = x(\theta)$, $y = y(\theta)$)

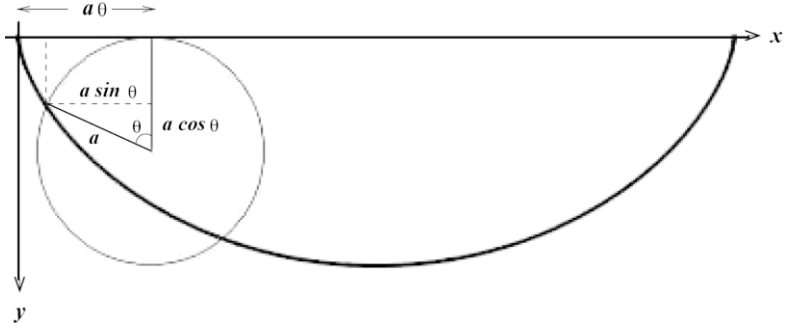


Fig. 7.1: The cycloid represented by the parametric equations in Eq. (7.1) with the y -axis pointing downwards. The geometrical interpretation of the parametric form in Eq. (7.1) is obvious from the figure. Note that at the extremities like, for e.g., near the origin, $y \propto x^{2/3}$ and hence the slope diverges.

to a cycloid has the form

$$x = a(\theta - \sin \theta); \quad y = a(1 - \cos \theta), \quad (7.1)$$

where a is the radius of the rolling circle. We shall now take a closer look at this result.

No doubt, there is progress

While the initial solution to the brachistochrone problem engaged some of the intellectual giants of seventeenth century, it is now within the grasp of an undergraduate student. Let $y(x)$ denote the equation to the curve which is the solution to the brachistochrone problem with the coordinates chosen such that x is horizontal and y is measured vertically downwards as in Figure 7.1. Let the particle begin its slide from the origin with zero velocity. If the infinitesimal arc length along the curve around the point $P(x, y)$ is $ds = (1 + y'^2)^{1/2} dx$ where $y' = (dy/dx)$, then the particle takes a time $dt = ds/v$, where $v = \sqrt{2gy}$ is its speed at P . To determine the curve we only need to find the extremum of the integral over dt , which is a straightforward problem in the calculus of variations. (In fact, if we replace earth's gravity by some other potential field, we only have to replace $v = \sqrt{2gy}$ by $v = [2(V_0 - V)/m]^{1/2}$.) Let us, however, analyse it from two slightly different approaches [24].

Trick 1: Use clever coordinates

In the first approach, we will make a coordinate transformation which simplifies the problem considerably. Let us introduce, in the first quadrant, two new coordinates α and β in place of the standard Cartesian coordinates (x, y) by the relations

$$x = \alpha^2 \left(\frac{\beta}{\alpha} - \sin \frac{\beta}{\alpha} \right); \quad y = \alpha^2 \left(1 - \cos \frac{\beta}{\alpha} \right), \quad (7.2)$$

where $\alpha > 0$ and $0 \leq \beta \leq 2\pi\alpha$. Obviously, for a fixed α , the curve $x(\beta), y(\beta)$ is a cycloid (which tells you that we are cheating a little bit using our knowledge of the final solution!). The square of the velocity of the particle

$$v^2 = 2gy = \dot{x}^2 + \dot{y}^2, \quad (7.3)$$

where overdots denote differentiation with respect to time, can now be expressed in terms of β and $\dot{\alpha}$ by straightforward algebra. This gives the relation

$$2gy = 2y\dot{\beta}^2 + 4 \left(2\alpha \sin \frac{\beta}{2\alpha} - \beta \cos \frac{\beta}{2\alpha} \right)^2 \dot{\alpha}^2. \quad (7.4)$$

The term involving $\dot{\alpha}^2$ is non-negative; further, since $y > 0$ we have $\dot{\beta} \leq \sqrt{g}$. Integrating this relation between $t = 0$ and $t = T$ where T is the time of descent, we get

$$\beta(T) = \int_0^T \dot{\beta} dt \leq \int_0^T \sqrt{g} dt = \sqrt{g}T. \quad (7.5)$$

It follows that the time of descent is bounded from below by the equality $T \geq \beta(T)/\sqrt{g}$. The best we can do is to set $\dot{\beta} = \sqrt{g}$ and $\dot{\alpha} = 0$ to satisfy Eq. (7.4) and hit the lower bound in Eq. (7.5). Since the required curve has $\alpha = \text{constant}$, it is obviously a cycloid parameterized by β .

The angular parameter of the cycloid, $\theta = \beta/\alpha$, varies with time at a constant rate $\dot{\theta} = \dot{\beta}/\alpha = \sqrt{g}/\alpha$. It is clear from the parameterization in Eq. (7.2) that the radius a of the circle which rotates to generate the cycloid is related to α by $a = \alpha^2$. Hence the angular velocity of the rolling circle is $\omega = \dot{\theta} = \sqrt{g}/a$. If the particle moves all the way to the other end of the cycloid at a horizontal distance $L = 2\pi a$, then the time of flight will be $T = 2\pi/\omega = (2\pi L/g)^{1/2}$. If L is 100 m, then with $g = 9.8 \text{ m s}^{-2}$ we get $T \approx 8 \text{ sec}$ which is better than the world record for a 100 m dash! Gravity seems to do quite well.

The 100 meter dash by gravity

Another indirect way of arriving at the cycloidal solution is also of some interest. This approach uses the concept of the hodograph which is the curve traced by a particle in the velocity space (see Chapter 3). Let us try to determine the hodograph corresponding to the motion of swiftest descent. For simplicity, consider the full transit of the particle from a point A to a point B in the same horizontal axis $y = 0$. Let the speed of the particle be v when the velocity vector makes an angle θ with respect to the v_x -axis in the velocity space. Then the hodograph is given by some curve $u(\theta)$ which we are trying to determine. Using $\dot{x} = v \cos \theta, \dot{y} = v \sin \theta$, $y = v^2/2g$, we can write the relations:

*Trick 2:
Use hodograph*

$$dt = \frac{dv}{g \sin \theta}; \quad dx = \frac{v dv}{g} \cot \theta. \quad (7.6)$$

We are now required to minimize the integral over dt while keeping the integral over dx fixed. Incorporating the latter constraint by a Lagrange multiplier $(-\lambda)$, we see that we need to minimize the following integral:

$$I = \int \frac{dv}{g} \left(\frac{1}{\sin \theta} - \lambda v \cot \theta \right). \quad (7.7)$$

The minimization is trivial since no derivatives of the functions are involved and leads to the relation $v = (1/\lambda) \cos \theta$ with $-\pi/2 < \theta < \pi/2$. We can now trade off the Lagrange multiplier λ for the total horizontal distance L (obtained by integrating dx) and obtain $\lambda^2 = \pi/2gL$. Hence, our hodograph is given by the equation

$$v(\theta) = \sqrt{\frac{2gL}{\pi}} \cos \theta \equiv 2R_0 \cos \theta. \quad (7.8)$$

This is just the polar equation for a circle of radius R_0 with the origin coinciding with the left-most point of the circle. (We saw earlier in Chapter 3 that the hodograph for the Kepler problem is also a circle but that was for motion in a $(1/r)$ potential; here we are studying the motion under the action of a constant gravitational field.)

How can we get to the curve in real space from the hodograph in the velocity space? In this particular case, it is quite easy. Suppose we shift the circular hodograph horizontally to the left by a distance R_0 . This requires subtracting a horizontal velocity which is numerically equal to the radius of the hodograph. After the shift, we obtain the hodograph of uniform circular motion, which is, of course, a circular hodograph with the origin at its center. Hence, the motion that minimizes the time of descent is just uniform circular motion added to a rectilinear uniform motion with a velocity equal to that of circular motion. This is, of course, the path traced by a point on a circle that rolls on a horizontal surface which is a cycloid. The advantage of this approach is that we obtain the cycloid in terms of its geometrical definition, instead of its equations.

From there to here

Box 7.1: Brachisto and other chrones: A bit of history

This tautochrone problem has appeared in English literature! Herman Melville's 1851 classic *Moby Dick* has a chapter called "The Try-Works" which describes how the try-pots of the ship *Pequod* are cleaned. (In case you haven't read the book, a try-pot is a large cauldron, usually made of iron, which is used to obtain liquid oil from whale blubber.) In that occurs the passage: "It was in the left hand try-pot of the *Pequod* that I was first indirectly struck by the remarkable fact, that in geometry all bodies gliding along the cycloid, my soapstone for example, will descend from any point in precisely the

Moral: Read Classics!

same time.” The remarkable fact Melville writes about is, of course, the tautochrone problem.

One of the early investigations about the time of descent along a curve was by Galileo. He, like many others, was interested in the time taken by a particle to perform an oscillation on a circular track which, of course, is what a simple pendulum of length L hanging from the ceiling will do. Today we could write down this period of oscillation as

$$T = \sqrt{\frac{L}{g}} \int_0^{\pi/2} \frac{d\theta}{\sqrt{1 - k^2 \sin^2 \theta}}, \quad (7.9)$$

where k is related to the angular amplitude of the swing. Of course, in the days before calculus, the expression would not have meant anything! Instead, Galileo used an ingenious geometrical argument and — in fact — thought that he had proved the circle to be the curve of fastest descent. It was, however, known to mathematicians of the 17th century that Galileo’s argument did not establish such a result.

The major development as regards the brachistochrone came up when Bernoulli threw a challenge in 1697 to the mathematicians of that day with the announcement I quoted in the beginning of this chapter. Bernoulli, of course, knew the answer and the problem was also solved by his brother Jakob Bernoulli, Leibniz, Newton and L’Hospital. Newton is said to have received Bernoulli’s challenge at the Royal Society of London one afternoon and (according to second hand sources, like John Conduit — the husband of Newton’s niece), Newton solved the problem by night-fall. The “solution”, which was simply a description of how to construct the relevant cycloid, was published anonymously in the *Philosophical Transactions of the Royal Society* of January 1697 (back dated by the editor Edmund Halley). Newton actually read aloud his solution in a Royal Society meeting only on 24 February 1697. Legend has it that Bernoulli immediately recognized Newton’s style and exclaimed “*tanquam ex ungue leonem*” meaning “the lion is known by its claw”.

The fact that brachistochrone and the tautochrone problems lead to the same curve, viz., the cycloid, in the case of a constant gravitational field is a bit of an accident. In general, if the potential varies as the square of the arc length along a curve, then a bead sliding on that curve will oscillate with a period independent of the amplitude. If the force field is constant, so that the potential is linear in the height, then this condition translates to a curve whose height should be proportional to the square of the arc length. It is straightforward to show that this condition is satisfied by the cycloid. In this sense, the tautochrone problem is rather trivial and only involves the force

*Physics of brachisto
and tauto are quite
different*

acting *along the curve* and is independent of the force *acting normal to the curve*. The situation regarding the brachistochrone is more complex. In this case, there should be a delicate balance between the centripetal force at any given point in the curve and the component of the external force perpendicular to the curve.

You should also bear in mind the following distinction when you think of the cycloid as a solution to both the tautochrone and brachistochrone motions. Given a cycloid, if you start a particle sliding from rest from any point, it will, of course, oscillate with a period independent of amplitude. But a particle starting at some arbitrary point in the cycloid will not be the correct extremal path for the brachistochrone problem. The correct cycloid that is the solution to the brachistochrone problem, for a particle starting from rest, always has the cycloid kink at the starting position.

Another difference

*Curves of
complementary
descent, defined*

Given the solution to the brachistochrone problem, one is naturally led to ask the following question: Let us consider a particle sliding along a given curve from the origin to a point (r, θ) taking the time $T(r, \theta)$. We want to know whether there exists another curve connecting these two points, on which the particle can slide, taking the same amount of time. Obviously, unless the first curve is a cycloid connecting the two points, we will expect to find alternative solutions. Such curves are called complementary curves of descent [25]. If $\theta = \theta(r)$ is the equation to the curve, then the time of descent is given by the integral of ds/v where s is the arc length and $v = \sqrt{2gy} = \sqrt{gr \sin \theta}$ is the velocity. Equating this to the given time of descent $T(r, \theta)$ we get the equation

$$\int \sqrt{\frac{1 + r^2 (d\theta/dr)^2}{2gr \sin \theta}} dr = T(r, \theta) . \quad (7.10)$$

Differentiating both sides with respect to r and manipulating the terms lead to a quadratic equation

$$A \left(\frac{d\theta}{dr} \right)^2 + B \frac{d\theta}{dr} + C = 0 , \quad (7.11)$$

with

$$\begin{aligned} A &\equiv 2gr \sin \theta \left(\frac{\partial T}{\partial \theta} \right)^2 - r^2 , \\ B &\equiv 4gr \sin \theta \frac{\partial T}{\partial r} \frac{\partial T}{\partial \theta} , \\ C &\equiv 2gr \sin \theta \left(\frac{\partial T}{\partial r} \right)^2 - 1 . \end{aligned} \quad (7.12)$$

This allows you to figure out complementary curves of descent of different kinds.

As a simple example, let the original curve be a straight line which makes an angle θ with respect to the x -axis. The time of descent in this case is given by the function

The strange complement to a straight line

$$T(r, \theta) = \sqrt{\frac{2r}{g \sin \theta}} . \quad (7.13)$$

We want to find a curve which is the complement to this, having the same time of descent. If you solve Eqs. (7.11), (7.12) with this function, you find that the solution is given by

$$r = 2b\sqrt{\cos \theta \sin \theta} , \quad (7.14)$$

which goes by the name Lemniscate of Bernoulli. Unfortunately this does not have any other interesting applications in physics.

There is a nice generalization of the brachistochrone problem which has not received much attention. The cycloid solution was obtained under the assumption of a uniform, constant gravitational field of a flat Earth. In reality, of course, the gravitational field varies as $(1/r^2)$ around a spherical object. The question arises as to how the curve of swiftest descent gets modified when we work with the $(1/r^2)$ force.

The brachistochrone for the $1/r^2$ field

To tackle this problem, it is convenient to use the polar coordinates in the plane of motion and approximate the gravitational source as a point particle of mass M at the origin. We are interested in determining the curve $r(\theta)$ such that a particle starting from a point A (with coordinates $r = R$ and $\theta = 0$) will reach a point B (with coordinates $r = r_f$, $\theta = \theta_f$) in the shortest possible time. We will, as usual, encounter some curious features.

The mathematical formulation of the variational principle is quite simple. If $v(r)$ is the speed of the particle when it is at the radial distance r , then

Maths is routine

$$v^2 = 2GM \left(\frac{1}{r} - \frac{1}{R} \right) = C^2 \left(\frac{1}{x} - 1 \right) , \quad (7.15)$$

where $x = r/R$ and $C^2 = 2GM/R$. The variational principle requires us to minimize the integral over ds/v where $ds = R d\theta (x'^2 + x^2)^{1/2}$ is the arc length along the curve with $x' = dx/d\theta$. This, in turn, requires determining the extremum of the integral

$$T = \frac{R}{C} \int d\theta \left(\frac{x'^2 + x^2}{(1/x) - 1} \right)^{1/2} \equiv \int L(x', x) d\theta . \quad (7.16)$$

The Euler-Lagrange equation will lead to a second order differential equation involving $x''(\theta)$. But since the integrand is independent of θ (“time”), we know that $x'(\partial L / \partial x') - L$ is conserved (“energy”). Equating it to a

constant K gives a first integral thereby allowing the problem to be reduced to quadrature. Fairly straightforward algebra then leads to the form of the function $\theta(x)$ given by the integral

$$\theta(x) = \int_1^x \frac{dy}{y} \sqrt{\frac{1-y}{\lambda y^3 + y - 1}}, \quad (7.17)$$

where $\lambda \equiv (R/KC)$.

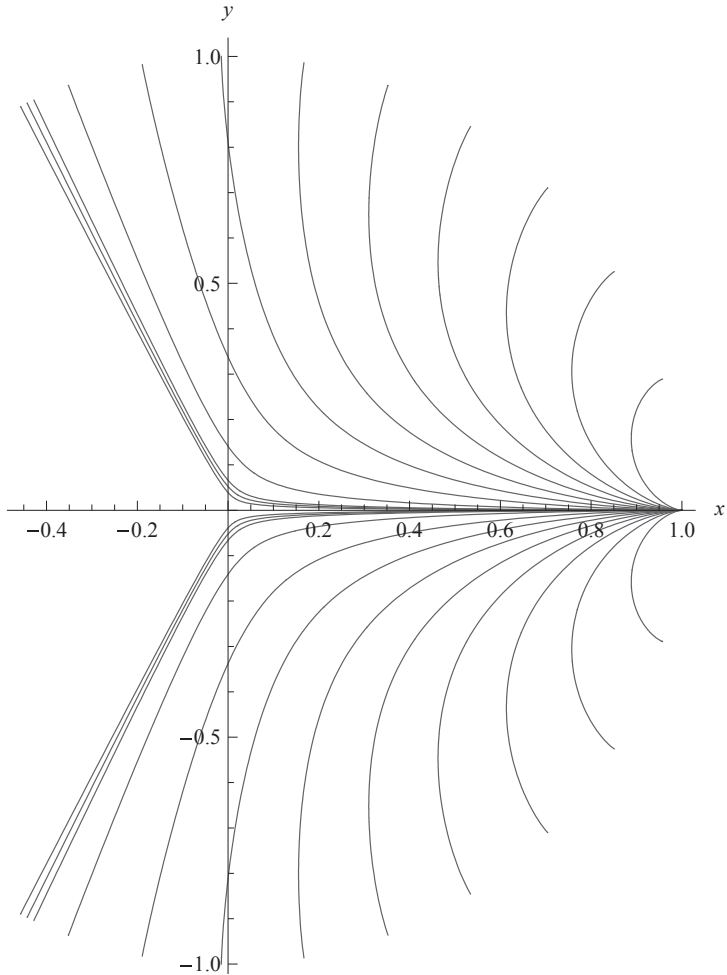


Fig. 7.2: Each of the curves in the figure gives the solution to the brachistochrone problem when the gravitational force falls as $(1/r^2)$ from the origin. Note that each curve has a turning point and none of the curves go through the “forbidden” region between $\theta = -(2\pi/3)$ and $\theta = +(2\pi/3)$.

Unfortunately, this is an elliptic integral making further analytic progress difficult. Working things out numerically, one can plot the relevant curves which show a very interesting pattern (see Figure 7.2). To begin with, one notices that each curve has a turning point $x = \ell$, say, where $(dx/d\theta) = 0$. This is a point of minimum approach related to λ by $\lambda = \ell^{-3}(1 - \ell)$. What is curious is the asymptotic behaviour of the curve after it turns around. It is clear from Figure 7.2 that the curves never enter the “forbidden region” between $\theta = -2\pi/3$ and $\theta = +2\pi/3$. This is obvious from the figure; but can we understand this analytically?

Surprise: The forbidden zone

One can do this but it requires a rather careful handling [26, 27] of the integral in Eq. (7.17). As you can easily see, what we need to prove is that the limiting value of θ given by this integral, when $\ell \rightarrow 0$ reaches a finite limit. To do this, let us rewrite Eq. (7.17) for $x = \ell$ after expressing λ as $\ell^{-3}(1 - \ell)$. This will give

Handle with care

$$\begin{aligned}\theta(\ell) &= \int_1^\ell \frac{dy}{y} \left[\frac{1-y}{\ell^{-3}(1-\ell)y^3 + y - 1} \right]^{1/2} \\ &= \left(\frac{\ell^3}{1-\ell} \right)^{1/2} \int_1^\ell \frac{dy}{y} \left[\frac{1-y}{(y-\ell)(y^2 + \ell y + \ell^2(1-\ell)^{-1})} \right]^{1/2}. \quad (7.18)\end{aligned}$$

The second relation is obtained by factorizing the denominator since $(y - \ell)$ must be one of its roots. We are interested in the $\ell \rightarrow 0$ limit of this integral which requires one more rescaling. Substituting $q = (\ell/y)^{3/2}$, our integral can be transformed to the form

$$\theta(\ell) = \frac{2}{3\sqrt{1-\ell}} \int_1^{\ell^{2/3}} dq \left\{ \frac{1 - \ell q^{-2/3}}{q(q^{1/3} - q)(q^{-4/3} + q^{-2/3} + (1-\ell)^{-1})} \right\}^{1/2}. \quad (7.19)$$

This one has a simple limit when $\ell \rightarrow 0$ and we get

$$\theta(0) = \frac{2}{3} \int_1^0 \frac{dq}{(1-q^2)^{1/2}} = -\frac{\pi}{3}. \quad (7.20)$$

The angle from the positive x -axis is $(\pi - \pi/3) = 2\pi/3$ because we have considered only the branch from the turning point. Further, there is a mirror symmetric curve in the lower half plane as well. So we find that when $\ell \rightarrow 0$ the angle of the trajectory reaches the asymptotic values:

$$\theta_{\text{crit}} = \mp \frac{2\pi}{3}. \quad (7.21)$$

In fact, the 3 in $(2\pi/3)$ of the forbidden zone comes from the power law index of the force. For the brachistochrone problem in r^{-n} force law, the forbidden zone is given by $-2\pi/(n+1) < \theta < 2\pi/(n+1)$.

*A nice problem,
with no name!*

Having described the classic variational problem which started it all, we now discuss another one, which does not even seem to have a respectable name. This problem [18] can be stated as follows. Consider a planet of a given mass M and volume V and a constant density $\rho = M/V$. We are asked to vary the shape of the planet so as to make the gravitational force exerted by the planet on a given point at its surface as high as possible. What is the resulting shape?

Most people would guess that the shape is either a sphere or something like the apex of a cone. The latter is easy to refute since it puts a fair amount of the mass away from the chosen point; but a sphere remains an intriguing possibility. The correct answer, however, is quite strange and can be obtained as follows.

Let the chosen point be at the origin and let the z -axis be along the direction of the maximal force acting on a test particle at the origin. It is obvious that this z -axis must be an axis of symmetry for the planet; if it is not, then one can increase the z -component of the net force by moving material from larger to smaller transverse distance until the planet is axially symmetric. So, our problem reduces to determining the curve $x = x(z)$ (with $0 < z < z_0$, say) which, on revolution around the z -axis, generates the surface of the planet. (The solution is plotted as a thick unbroken curve in [Figure 7.3](#).)

To calculate the z -component of the gravitational force acting on the origin, we divide the planet into circular discs, each of thickness dz , lo-

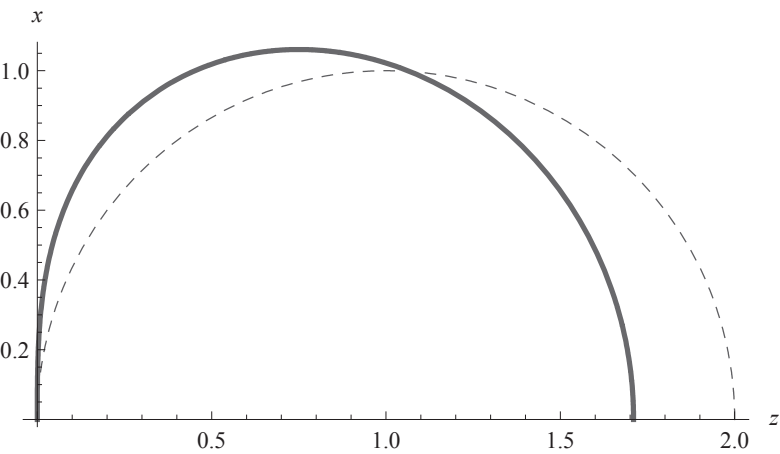


Fig. 7.3: The solid of revolution obtained by rotating the unbroken (thick) curve, about the z -axis, will give the shape of a constant density planet that will exert the maximum possible z -component of gravitational force at the origin. This shape does not seem to have any special name. The dashed (thin) curve is a sphere with the same volume given for comparison.

cated perpendicular to z -axis. To get the force exerted on a test particle of mass m by any single disc, we further divide it into annular rings of inner radii x and outer radii $x + dx$. The force along the z -axis by any one such ring will be given by

Routine maths but ...

$$dF = Gm(\rho 2\pi x dx) dz \frac{1}{x^2 + z^2} \frac{z}{\sqrt{x^2 + z^2}} . \quad (7.22)$$

Hence the total force is given by

$$\begin{aligned} F &= 2\pi Gm\rho \int_0^{z_0} dz \int_0^{x(z)} x dx \frac{z}{(x^2 + z^2)^{3/2}} \\ &= \frac{3GMm}{2a^3} \int_0^{z_0} dz \left(1 - \frac{z}{(x^2(z) + z^2)^{1/2}} \right) . \end{aligned} \quad (7.23)$$

In arriving at the last expression we have expressed the density as $\rho = 3M/4\pi a^3$ so that the volume of the planet is constrained by the condition

$$V = \pi \int_0^{z_0} dz x^2(z) = \frac{4\pi a^3}{3} . \quad (7.24)$$

Imposing this condition by a Lagrange multiplier $(-\lambda)$, we see that we have to essentially find the extremum of the integral over the function

$$L = 1 - \frac{z}{(x^2 + z^2)^{1/2}} - \lambda x^2 . \quad (7.25)$$

This is straightforward and we get

$$\frac{z}{(x^2 + z^2)^{3/2}} = 2\lambda = \frac{1}{z_0^2} , \quad (7.26)$$

where the last equality determining the Lagrange multiplier follows from the condition that when $z = z_0$ we have $x = 0$. Our constraint on the total volume [given by Eq. (7.24)] implies that $z_0^3 = 5a^3$ thereby completely solving the problem. The polar equation to the curve is

$$r^2 = 5^{2/3} a^2 \cos \theta ; \quad (7.27)$$

for comparison, a sphere with the same volume will be described by the equation $r = 2a \cos \theta$.

With hindsight, one can obtain this result from a simpler, intuitive argument. The crucial point is to realize that all the small elements of mass dm on the surface of the material must contribute *equally* to the z -component of the force at the origin. If this is not the case, we can simply move a small amount of matter from one point to another point on the surface thereby increasing the force. If we denote the mass element by their distance r from the origin and the angle θ which it makes with the z -axis, then an in-

... you could have got it with no maths!

infinitesimal element of mass dm on the surface provides the z -component of the force which varies as $F_z = (Gdm/r^2)\cos\theta$. Since this has to be independent of the location, the surface must satisfy $r^2 \propto \cos\theta$ which is precisely our solution.

The shape of our weird planet is shown in Figure 7.3 by the thick unbroken curve (along with that of a sphere with same volume) which has no cusps at the poles. This shape does not seem to have any specific name. The total force exerted by this planet at the origin can be computed using Eq. (7.23). We get:

$$F = \left(\frac{27}{25}\right)^{1/3} \frac{GMm}{a^2} \approx 1.03 \frac{GMm}{a^2}, \quad (7.28)$$

But then, it is the principle that matters.

which is not too much of a gain over a sphere.

We note a minor subtlety which we glossed over while doing the variation in this problem. Unlike the usual variational problems, the end point z_0 is not given to us as fixed while doing the variation of the integrals in Eq. (7.23), Eq. (7.24). It is possible to take this into account by a slightly more sophisticated treatment but it will lead to the same result in this particular case.

A beauty from extremum

Another beautiful phenomena all of us are familiar with, which owes its existence to an extremum principle, is the rainbow. We all know that a rainbow is formed when the light from the Sun that is scattered by a raindrop reaches your eye. But, of course, there are raindrops all over the sky, while you see the rainbow at a characteristic angle and shape in the sky! This is due to the fact that you will see the rainbow only when a large number of rays of light are accumulating in a particular direction after passing through the raindrop.

Figure 7.4(a) shows the path of a light ray through a spherical droplet of water, which leads to the formation of, what is called, a primary rainbow. The ray incident at A gets refracted; part of the light is reflected at B which is again refracted at C. The angles x and y are related by $\sin x = n \sin y$ where n is the refractive index of water. The direction of the ray changes by $(x - y)$ at A, by $(\pi - 2y)$ at B and by $(x - y)$ at C thereby undergoing a total deviation $D(x) = 2x - 4y + \pi$.

The net effect of the water droplet is to deviate a ray of light as shown in Fig. 7.4(b), where the incident direction of ray is taken to be horizontal. The angle of incidence x will be different for droplets of water at different locations and, in general, D will change with x . There is, however, one particular angle x_c at which $(dD/dx) = 0$. At this critical value, the deviation $D = D_c$ is stationary with respect to x and one sees an enhancement of several rays traveling towards the same direction after going through the water droplets. (In the above analysis, we only maximize the deviation angle D with respect to the incident angle x . Rigorously speaking, we have to worry about the cross section of the raindrops available to the

The critical ray

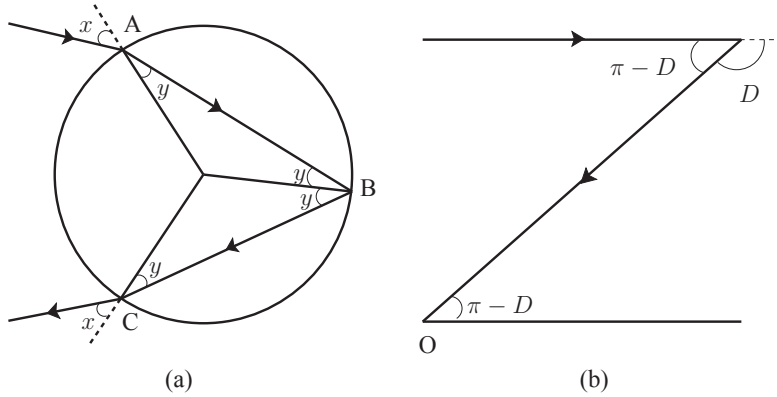


Fig. 7.4: (a) The path of the light ray through a raindrop which produces the primary rainbow. The net effect of two refractions (at A, C) and one reflection (at B) is to deviate the light ray by an angle $D = 2x - 4y + \pi$. (b) At a critical angle of incidence, D is an extremum with respect to x and a large number of rays accumulate along this direction undergoing a deviation D_c . This causes a rainbow in the sky located on the semicircular rim of a cone with vertex at O and semi vertical angle $(\pi - D_c) = 4y_c - 2x_c$.

light incident at different angles which, in the case of the spherical geometry, is governed by the usual $\sin \theta d\theta d\phi$ factor. Fortunately, this does not affect the final conclusion.) This will lead to a rainbow in the sky located on the semicircular rim of a cone with vertex at O and semi vertical angle $(\pi - D_c) = 4y_c - 2x_c$. Elementary calculation now gives

$$\cos^2 x_c = \frac{1}{3}(n^2 - 1). \quad (7.29)$$

Taking the refractive index for $\lambda = 400$ nm to be $n_{400} = 1.3440$ and for $\lambda = 700$ nm to be $n_{700} = 1.3309$, we find that $x_c = (58.77^\circ, 59.54^\circ)$ and $y_c = (39.51^\circ, 40.36^\circ)$ for the two wavelengths, leading to $(\pi - D_c) = (40.51^\circ, 42.38^\circ)$. Thus, the primary rainbow is at about 41° and its angular width is about 1.87° .

All that beauty, just from a few numbers

A little thought shows that while it is possible for a raindrop to scatter light at values smaller than 42° , it cannot do it at angles larger than 42° . This has the consequence that the region in the sky below the rainbow appears brighter than the region above it.

It is now obvious that one can obtain similar results with the light rays reflecting more than once inside the raindrop. This leads to what is known as secondary, tertiary etc. rainbows in the sky. It is easy to repeat the analysis in these cases and we will find that, for the N th order rainbow, Eq. (7.29) gets replaced by the result

Given 1, make 2, 3, ...

$$\cos^2 x_N = \frac{1}{N(N+2)}(n^2 - 1), \quad (7.30)$$

which, of course, reduces to Eq. (7.29) when $N = 1$.

For $N = 2$, we get the secondary rainbow at an angle of 52° which is about 10° higher in the sky than the primary. It is less bright (by about 43 per cent) than the primary because of the additional loss of intensity due to the second reflection. The second reflection also reverses the colour sequence in the secondary; the red edge of the rainbow will appear lower in the sky than the violet one.

The real surprise is with the tertiary

The geometry gets a bit trickier when we move to $N = 3$. The total deviation suffered by the light ray is now 318.4° after 3 reflections. This means that the tertiary rainbow is actually *behind* you — and is a circular halo around the Sun at about 41.6° — when you are facing the primary and secondary rainbows! If you proceed along these lines, the position of the first six orders of rainbows in the sky around you will be as shown in the Fig. 7.5.

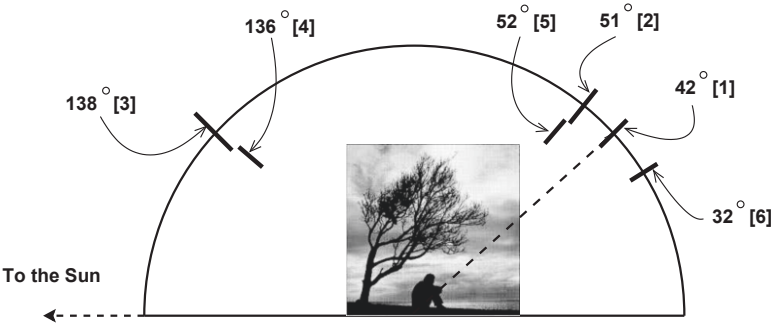


Fig. 7.5: The locations of different orders of rainbow in the sky marked in square brackets as [1], [2], ... etc. The primary rainbow is at around 42° and the secondary one is at 51° . The next two, the tertiary and the fourth order rainbows, are *behind* the observer when she is facing the primary rainbow! The fifth and sixth order rainbows are in the forward direction, but unfortunately too faint to be seen.

Box 7.2: Rainbow – through the ages

Given the rather spectacular visual nature of the rainbow and the fact that it is not a periodic phenomenon in the sky (unlike for e.g., the waxing and waning of the Moon or the orbits of celestial orbits), it is no surprise that it had attracted considerable attention from pre-history. (For a detailed description of the history, see Ref. [28].)

Many people have provided “explanations” for the rainbow, including Aristotle, Kepler and Gilbert. The clearest and the one closest to the correct explanation came in the middle ages, both in the east and in the west. In the east, it was due to Kamal al-Din al-Farisi (1267–1319) and in the west it was from the German monk Dietrich

von Freiberg (1250–1311). Both of them correctly stated that the scattering is due to individual raindrops — unlike many before them who thought it was from the rain clouds themselves. Freiberg was also the first to associate the primary rainbow with two refractions and one reflection and the secondary rainbow with two refractions and two reflections in his small book, *De iride et radialidus impressionibus*.

The only thing missing in these explanations is the fact that rays get concentrated at a particular degree. This was the major contribution from Descartes and he obtained this by actually tracing the rays through spherical water droplets with pencil, paper and, of course, Snell's law. This way he obtained both the primary and secondary rainbows and their respective angles of 42° and 52° .

The secondary rainbow, as far as I know, has been described in the contemporary literature only once and at that time the author got it wrong! Rebecca Goldstein ends her novel "Strange attractors" [29] describing a group of mathematicians going outdoors to look at a double rainbow. She puts the secondary rainbow "beneath" the primary one, though with the correct inversion of the spectrum.

All these naturally suggested the existence of higher order rainbows and many intrigued people searched the sky in vain for centuries, particularly for the tertiary rainbow. Being quite logical, they were all looking in the sky above the secondary, maybe another 10° up. For reasons which will be obvious to you from the previous discussion, nobody ever saw it in the historical days.

It is unclear whether even Newton, who worked out all the details of the n th order rainbow, bothered to actually calculate the specific angular position of the tertiary rainbow; if he did, he did not publish it either in a series of inaugural lectures as Lucasian Professor in 1670-72 or in his work *Opticks*. In the latter, he merely says that the light that undergoes three or more reflections is "scarcely strong enough to cause a sensible bow". Of course we know that — since the tertiary rainbow is a halo around the Sun — the glare of the Sun will completely wipe out this rainbow, making Newton's comment rather irrelevant if he had calculated the exact position. (Bernoulli also discusses this issue without identifying its location in the sky.) The clear statement as to where the tertiary is located and why it is impossible to see seems to have been first published by Halley as late as in the 1700s.

Obviously, you can hope to spot the tertiary rainbow only in a happy circumstance in which the Sun's glare is blocked. Eclipses are obvious choices but you also need to have rain as well as the proper angle for the sunlight. Given all these, it is not surprising that photographing the tertiary rainbow was not achieved until as late as 2011!

*Hard work without
calculus*

*Fiction is, after all,
fiction*

*... in spite of some
occasional claims to
the contrary!*

*Did Newton know
where to look?*

Got it, at last!

Michael Grossmann [30] happened to witness a rain shower in south-west Germany on 15 May 2011. The rain was falling sun ward while a dark cloud and a tree blocked part of the intensity in the sky near the Sun. On that day, Grossman managed to get a photograph of the tertiary rainbow which was in agreement with the theory!

The simplest form of the fluid flow, that arises when a body moves through a hypothetical fluid, will satisfy the following conditions: First, the fluid is assumed to be incompressible with the density being a constant. Then, the conservation of mass, expressed in the form of the continuity equation

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{v}) = 0, \quad (8.1)$$

(in which ρ is the density and \mathbf{v} is the fluid velocity) reduces to the simple condition $\nabla \cdot \mathbf{v} = 0$. Second, we will assume that the flow is irrotational ($\nabla \times \mathbf{v} = 0$) allowing for the velocity to be expressed as a gradient of a scalar potential $\mathbf{v} = \nabla \phi$. Finally we will ignore all properties of real fluids, like viscosity, surface tension etc. and will treat the problem as one of finding the solutions to the two equations $\nabla \cdot \mathbf{v} = 0$ and $\nabla \times \mathbf{v} = 0$ subject to certain boundary conditions. Equivalently, we find that the potential satisfies Laplace's equation $\nabla^2 \phi = 0$. So, the problem reduces to solving the Laplace equation with \mathbf{v} satisfying the boundary conditions — which are the only non-trivial features of the problem! Such a problem is considered to be well-understood but — as we will see in this chapter — even the simplest of them can lead to surprises [31].

*Flow of dry water =
Electrostatics*

Let us consider a body of an arbitrary shape moving through the fluid with a velocity \mathbf{u} . Then we need to solve the Laplace equation subject to the boundary condition $\mathbf{n} \cdot \mathbf{v} = \mathbf{n} \cdot \mathbf{u}$ at the surface, where \mathbf{n} is the normal to the surface. We would expect the fluid flow near the body to be affected by its motion but this effect should be negligible at sufficiently large distances. Hence the fluid velocity \mathbf{v} will be zero at spatial infinity.

The general form of the fluid velocity at large distances from the body (of arbitrary shape) can be determined by the following argument. We know that the function $1/r$ satisfies the Laplace equation. Further, if ϕ satisfies the Laplace equation, the spatial derivatives of ϕ also satisfy the same equation. Therefore, the directional derivative of $1/r$, along some

A cute result

direction specified by an arbitrary vector \mathbf{A} will also satisfy the Laplace equation. Such a directional derivative is given by $\mathbf{A} \cdot \nabla(1/r)$ and will fall as $1/r^2$ at large distances. Hence, at large distances from the body, we can take the leading order terms in the potential to be

$$\phi = -\frac{q}{r} + \mathbf{A} \cdot \nabla \left(\frac{1}{r} \right) + \mathcal{O}(1/r^3) . \quad (8.2)$$

*A trick to get
dipole potential*

This, of course, can be recognized as just electrostatics in disguise; the expansion in Eq. (8.2) is just the large distance expansion of the potential due to a distribution of charges. The first term is the monopole Coulomb term and the second one is the dipole term. (Incidentally, the dipole term is just the difference in the potential due to two charges kept separated by a distance \mathbf{A} ; clearly, the net potential will be the directional derivative along \mathbf{A} . This is the quickest way to get the dipole potential.) At sufficiently large distances we can ignore further terms, obtained by taking the second, third, derivatives of $1/r$.

The velocity field is then the analogue of the electric field in electrostatics. From the Gauss law we know that the flux of the electric field at large distances is proportional to the ‘total charge’ q . At large distances, the flux of the velocity field in our problem vanishes. Hence, it follows that $q = 0$ and the asymptotic form of the potential must have the form:

$$\phi = \mathbf{A} \cdot \nabla \left(\frac{1}{r} \right) = -\frac{\mathbf{A} \cdot \mathbf{n}}{r^2} , \quad (8.3)$$

where \mathbf{n} is the unit vector in radial direction. Taking the gradient, we get the velocity field to be

$$\mathbf{v} = (\mathbf{A} \cdot \nabla) \nabla \left(\frac{1}{r} \right) = \frac{3(\mathbf{A} \cdot \mathbf{n})\mathbf{n} - \mathbf{A}}{r^3} . \quad (8.4)$$

Electrostatic insight

(These manipulation are most efficiently done using index notation and summation convention, with $\partial_\alpha r = (1/2r)\partial_\alpha r^2 = x^\alpha/r$ used repeatedly.) The actual form of \mathbf{A} needs to be determined using the conditions near the body (which will be a mess for a body of arbitrary shape) but it is interesting that the flow at large distances is fixed entirely in terms of a single vector \mathbf{A} . In fluid mechanics, it is a bit of a surprise but in electrostatics it is not. If the monopole vanishes, you would expect the dipole moment to determine the behaviour of electric field at large distances.

The *real* surprise comes up when we try to calculate the *total* kinetic energy associated with the fluid flow given by

$$K_{\text{lab}} = \frac{1}{2} \rho \int d^3\mathbf{x} v^2 , \quad (8.5)$$

where the integral is over all space outside a sphere of radius a and the subscript “lab” stands for the lab frame in which the sphere is moving with a velocity \mathbf{u} . (The fact that the sphere is moving is irrelevant since it only shifts the origin by \mathbf{ut} which is a constant as far as the spatial integration is concerned.) While the fluid flow at large distances can be expressed entirely in terms of a single vector \mathbf{A} , the flow closer to the body can be extremely complicated. Hence, one might have thought that, in such a general case, one cannot infer anything about the total kinetic energy of the fluid. But it is indeed possible to express the total kinetic energy of the fluid flow entirely in terms of the single vector \mathbf{A} even though the fluid flow everywhere *cannot* be expressed in terms of \mathbf{A} alone. (This result, as well as Eq. (8.8) and Eq. (8.20) below, are derived in Ref. [32] but do not seem to be discussed in detail in any other book.)

What is the total kinetic energy?

To obtain this result, we will use the identity $v^2 = u^2 + (\mathbf{v} + \mathbf{u}) \cdot (\mathbf{v} - \mathbf{u})$. If we integrate both sides of this equation over a large volume V , the first term on the right will give a contribution proportional to $(V - V_0)$, where V_0 is the volume of the body. In the second term, we write $(\mathbf{v} + \mathbf{u}) = \nabla(\phi + \mathbf{u} \cdot \mathbf{r})$. Using $\nabla \cdot \mathbf{v} = 0$, $\nabla \cdot \mathbf{u} = 0$, we can write the second term as a total divergence $\nabla \cdot [(\phi + \mathbf{u} \cdot \mathbf{r})(\mathbf{v} - \mathbf{u})]$. On integrating this over the whole space, the second term becomes a surface integral over the surface of the body *and* a surface at large distance. That is, we have proved:

$$\int v^2 dV = u^2(V - V_0) + \oint_{S+S_0} (\phi + \mathbf{u} \cdot \mathbf{r})(\mathbf{v} - \mathbf{u}) \cdot \mathbf{n} dS, \quad (8.6)$$

where S is a surface bounding the volume V at large distance and S_0 is the surface of the body and the surface integral is taken over both.

The $(\mathbf{v} - \mathbf{u}) \cdot \mathbf{n}$ term vanishes on the surface of the body, due to the boundary conditions; hence we get no contributions from there! This is good since we have no clue about the pattern of velocity flow near the body. On the surface at large distances from the body, we can use the asymptotic form of the velocity field given in Eq. (8.4) to perform the integral, taking the surface to be a sphere of large radius R . The area $dS = R^2 d\Omega$ increases as R^2 while v falls as $1/R^3$ and ϕ falls as $1/R^2$. So $\phi(\mathbf{v} - \mathbf{u}) \cdot \mathbf{n} \approx -\phi \mathbf{u} \cdot \mathbf{n}$ on S . Hence the surface integral in Eq. (8.6) on S becomes the sum

$$-\oint_S \phi \mathbf{u} \cdot \mathbf{n} R^2 d\Omega + \oint_S (\mathbf{u} \cdot \mathbf{n})(\mathbf{v} \cdot \mathbf{n}) R^3 d\Omega - \oint_S (\mathbf{u} \cdot \mathbf{n})^2 R^3 d\Omega. \quad (8.7)$$

The integration over angular coordinates can be done using the easily proved relation $\langle (\mathbf{A} \cdot \mathbf{n})(\mathbf{B} \cdot \mathbf{n}) \rangle = (1/3)\mathbf{A} \cdot \mathbf{B}$ where $\langle \dots \rangle$ denotes the angular *average* which is $1/4\pi$ times the integral over $d\Omega$. Using this, we see that the integral over $-(\mathbf{u} \cdot \mathbf{n})^2 R^3$ gives $-u^2 V$, which precisely cancels with the $u^2 V$ in the first term in Eq. (8.6). Using Eq. (8.3) and Eq. (8.4),

Kinetic energy is also fixed by \mathbf{A}

we get the final answer to be:

$$K_{lab} = \frac{1}{2} \rho (4\pi \mathbf{A} \cdot \mathbf{u} - V_0 u^2) . \quad (8.8)$$

Thus, if we know the motion of the fluid at very large distances from the body, we can compute the *total* kinetic energy of the fluid flow *without* ever knowing the velocity field close to the body!!

A curiosity

We can obtain another curious result using this. To do this, we note that the K_{lab} can also be expressed in a different form of surface integral. Writing $\mathbf{v} = \nabla \phi$ the expression for kinetic energy reduces to

$$K = \frac{1}{2} \rho \int_{\mathcal{V}} d^3 \mathbf{x} (\nabla \phi)^2 = \frac{1}{2} \rho \int_{\mathcal{V}} d^3 \mathbf{x} \nabla \cdot (\phi \nabla \phi) , \quad (8.9)$$

where we have used $\nabla^2 \phi = 0$. Using Gauss theorem, this expression can be converted to a surface integral over the body and over a surface at large distance. The second one vanishes, giving

$$K_{lab} = -\frac{1}{2} \rho \oint_{S_0} dS (\mathbf{n} \cdot \mathbf{v}) \phi = -\frac{1}{2} \rho \oint_{S_0} dS (\mathbf{n} \cdot \mathbf{u}) \phi , \quad (8.10)$$

where we have used $\mathbf{n} \cdot \mathbf{v} = \mathbf{n} \cdot \mathbf{u}$ at the surface. Using the expression for K_{lab} from Eq. (8.8), we can now obtain the following result for the integral of $(\mathbf{n} \cdot \mathbf{u}) \phi$ over the surface of the body:

$$-\oint_{S_0} dS (\mathbf{n} \cdot \mathbf{u}) \phi = (4\pi \mathbf{A} \cdot \mathbf{u} - V_0 u^2) , \quad (8.11)$$

even though we do not know either the shape of the body or the velocity potential on the surface!

Let us now look at the electrostatic analogue of this result. You are given a distribution of charges with $q_{tot} = 0$ and dipole moment \mathbf{p} in a region bounded by a surface S_0 . You are also given a constant vector \mathbf{E}_0 and you are told that the component of the electric field normal to S_0 is given by $\mathbf{n} \cdot \mathbf{E}_0$. Then, the electrostatic energy is proportional to $(4\pi \mathbf{p} \cdot \mathbf{E}_0 - V_0 E_0^2)$ where V_0 is the volume of the region bounded by S_0 .

*Simple case
of a sphere*

We will now specialize to the simplest of all possible shapes for the body: a sphere of radius a . In this case, the dipole potential happens to be the *exact* solution at all distances outside the sphere. This is not difficult to understand. Given the spherical symmetry, the only vector that can appear in the solution is the velocity of the body \mathbf{u} . Linearity of the Laplace equation (and the boundary condition) tells you that the potential must be linear in this vector \mathbf{u} . Hence the solution must have the form in Eq. (8.3) with $\mathbf{A} \propto \mathbf{u}$. Using the boundary condition $\mathbf{n} \cdot \mathbf{v} = \mathbf{n} \cdot \mathbf{u}$ at the surface, it is

easy to show that

$$\mathbf{A} = \frac{1}{2}a^3\mathbf{u}, \quad (8.12)$$

which completely solves the problem. We will now explore this solution.

Given the fluid flow pattern everywhere, we can explicitly compute the total kinetic energy carried by the flow using any of the expressions derived above. We get

Effective mass from kinetic energy

$$\begin{aligned} K_{\text{lab}} &= -\frac{1}{2}\rho \int a^2 d\Omega \left(-\frac{1}{a^2}\right) (\mathbf{A} \cdot \mathbf{n})(\mathbf{u} \cdot \mathbf{n}) \\ &= \frac{1}{2}\rho(4\pi)\frac{1}{3}(\mathbf{A} \cdot \mathbf{u}) = \frac{1}{4}m_{\text{disp}}u^2, \end{aligned} \quad (8.13)$$

where m_{dis} is the mass of the fluid displaced by the sphere. So the total kinetic energy is $(1/2)[m_{\text{body}} + (1/2)m_{\text{dis}}]u^2$, with the fluid adding $(1/2)m_{\text{dis}}$ to the effective mass of the sphere. Of course, our general expression, Eq. (8.8) leads to the same result when we use Eq. (8.12) and everything seems fine.

We next consider the total momentum \mathbf{P} carried by the fluid which is the integral over all space of $\rho\mathbf{v}$. Normally, we would have expected it to be $(1/2)m_{\text{disp}}\mathbf{u}$ but we are in for a rude shock. By symmetry, the vector \mathbf{P} has to be in the direction of \mathbf{u} so we only need to compute the scalar $\mathbf{P} \cdot \mathbf{u}$. But since v falls as $1/r^3$ and the volume grows as r^3 we are in trouble! (This did not happen for the kinetic energy since we were integrating $v^2 \propto 1/r^6$ over all space.) Explicitly, we have,

The misbehaving momentum

$$\begin{aligned} \mathbf{P}_{\text{lab}} \cdot \mathbf{u} &= \rho \int d^3\mathbf{x} \frac{1}{r^3} [3(\mathbf{A} \cdot \mathbf{n})(\mathbf{u} \cdot \mathbf{n}) - \mathbf{A} \cdot \mathbf{u}] \\ &= \rho \int_a^\infty \frac{dr}{r} \int d\Omega [3(\mathbf{A} \cdot \mathbf{n})(\mathbf{u} \cdot \mathbf{n}) - \mathbf{A} \cdot \mathbf{u}]. \end{aligned} \quad (8.14)$$

Obviously, our power counting argument is correct and the r -integral diverges logarithmically at large distances! On the other hand, the angular integration over spherical surfaces gives zero because $\langle 3(\mathbf{A} \cdot \mathbf{n})(\mathbf{u} \cdot \mathbf{n}) \rangle = \mathbf{A} \cdot \mathbf{u}$ cancels the second term. It is incredible that the simplest problem in fluid flow past a body actually leads to a product of zero and infinity!

Infinities? In fluid flow past a sphere?!

If we perform the integral between two spheres of radii $r = a$ and $r = R$ centered on the moving sphere at any given instant of time, then the answer is indeed zero because the angular average gives zero. This would have been an acceptable result, except for two reasons. First, the result depends on taking the outer boundary to be a sphere. If we choose some other shape, say, a cylinder coaxial with the direction of motion of the sphere, the result can be different. One feels uneasy about the result depending on what one is doing at infinity especially since the direction of \mathbf{u} breaks the spherical symmetry.

A way out, but a cheap one

Second, one can argue that, if the sphere is pushed (through a fluid) from rest until it acquires a velocity \mathbf{u} , then — in the process — some momentum is imparted to the fluid. To compute this, one needs to know the pressure which acts on the sphere when \mathbf{u} is a function of time [33]. Let me briefly indicate how this can be obtained. The starting point is the Euler equation

$$\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} = -\frac{\nabla p}{\rho}. \quad (8.15)$$

When $\mathbf{v} = \nabla \phi(t, \mathbf{x})$, you can manipulate this equation to show that

$$\nabla \left[p + \frac{1}{2} \rho v^2 + \rho \left(\frac{\partial \phi}{\partial t} \right) \right] = 0, \quad (8.16)$$

so that the pressure can be expressed in the form

$$p = p_\infty - \frac{1}{2} \rho v^2 - \rho \frac{\partial \phi}{\partial t}, \quad (8.17)$$

This is just the time dependent version of Bernoulli's equation.

where p_∞ is the pressure at infinity. We are interested in the net force in the direction of motion of the sphere, taken to be the z -axis, which can be obtained by integrating $p \cos \theta$ over the surface of the sphere. From Eq. (8.4) we see that v^2 will be a function of $\cos^2 \theta$ so the contribution from the first two terms in Eq. (8.17) will vanish on integration over a sphere. The only surviving contribution comes from the last term, which can be easily evaluated to give

$$F_z = - \int_0^\pi 2\pi a^2 \sin \theta d\theta \left[\frac{1}{2} \rho a \cos^2 \theta \frac{du_z}{dt} \right] = \frac{1}{2} m_{disp} \frac{du_z}{dt}. \quad (8.18)$$

Clearly, the total momentum imparted is

$$\int F_z dt = \frac{1}{2} m_{disp} u_z, \quad (8.19)$$

which makes sense when we remember that the kinetic energy comes with the effective mass $(1/2)m_{disp}$. So, this is another purely local reason to believe that the total momentum of the fluid flow is non-zero.

In fact, we can generalize this argument and obtain a finite expression for the momentum for *any* body moving through a fluid. This momentum, once again, can be expressed entirely in terms of the vector \mathbf{A} for a body of *arbitrary* shape. To obtain this result, we use Eq. (8.8) and the relation $dE = \mathbf{u} \cdot d\mathbf{P}$ which relates the infinitesimal changes in the energy and momentum. To prove this relation, let us assume that the body is accelerated by some external force \mathbf{F} causing the momentum of the fluid flow to increase by an amount $d\mathbf{P}$ in a time interval dt . From the relation $d\mathbf{P} = \mathbf{F} dt$, we immediately get $\mathbf{u} \cdot d\mathbf{P} = \mathbf{F} \cdot \mathbf{u} dt = dE$. Given the form of E , it is now

Another nice, general, result

an elementary matter to verify that the total momentum of the fluid flow is given by

$$\mathbf{P} = 4\pi\rho\mathbf{A} - \rho V_0\mathbf{u} . \quad (8.20)$$

We see that this is, in general, non-zero. In the case of the sphere it does give $(1/2)m_{dis}\mathbf{u}$ which what we naively would have expected. Of course, the argument is designed to give this.

When we study the same result in the rest frame of the sphere, it becomes more apparent that we need to regularize the problem by introducing a very large (but finite) volume for the total fluid. In this frame, we have a sphere of radius a located around the origin and the fluid is flowing past it. The boundary condition at infinity is now different and we expect the fluid velocity to reach a constant value $-\mathbf{u}$ at large distances. (In the electrostatic case, this is easily achieved by adding a constant electric field to a dipole.) This leads to a velocity potential of the form

Go to the rest frame of the sphere ...

$$\psi = -\mathbf{r} \cdot \mathbf{u} + \phi = -\mathbf{r} \cdot \mathbf{u} - \frac{\mathbf{A} \cdot \mathbf{n}}{r^2} . \quad (8.21)$$

We denote the velocity potential in the rest frame by ψ to distinguish it from the velocity potential in the lab frame, ϕ . Let us now ask what is the kinetic energy of the fluid in this frame in which the body is at rest. The fluid velocity now is $\mathbf{w} = \mathbf{v} - \mathbf{u}$. The kinetic energy in the rest frame will be

... and land in serious trouble again

$$\begin{aligned} K_{\text{rest}} &= \int d^3\mathbf{x} \frac{1}{2} \rho w^2 = \frac{1}{2} \rho \int d^3\mathbf{x} [v^2 + u^2 - 2\mathbf{v} \cdot \mathbf{u}] \\ &= \frac{1}{2} \rho \int d^3\mathbf{x} u^2 - \mathbf{u} \cdot \mathbf{P}_{\text{lab}} + K_{\text{lab}} . \end{aligned} \quad (8.22)$$

We see that the last term is the kinetic energy in the lab frame, K_{lab} , which is well-defined. The second term is ambiguous. It vanishes if we use spherical regularization, but is given by Eq. (8.20) if we use local energy conservation arguments. In the latter case, $K_{\text{lab}} - \mathbf{u} \cdot \mathbf{P}_{\text{lab}} = -(1/4)m_{disp}u^2$ is *negative*. The first term, however, will be divergent if we take the volume of the fluid to be infinite and is positive. This divergence arises because, if the fluid extends all the way to infinity, then most of it will be moving with a velocity $-\mathbf{u}$ in the rest frame of the sphere. This will contribute an infinite amount of kinetic energy. While quite understandable, it shows that Galilean invariance needs to be used with care in the presence of an external medium. There is no simple way of handling this difficulty.

Moral: Galilean invariance is tricky in a medium

I conclude this chapter with another, seemingly paradoxical, result in fluid flow which, fortunately, is well understood. But it leads to a curious, and not so well known effect. Consider the flow of a fluid through an orifice of area A_2 as shown in Fig. 8.1. We will assume that $A_1 \gg A_2$ and the fluid is incompressible giving $v_1 A_1 = v_2 A_2$ and hence $v_1 \ll v_2$. Using

Result from energy conservation ...

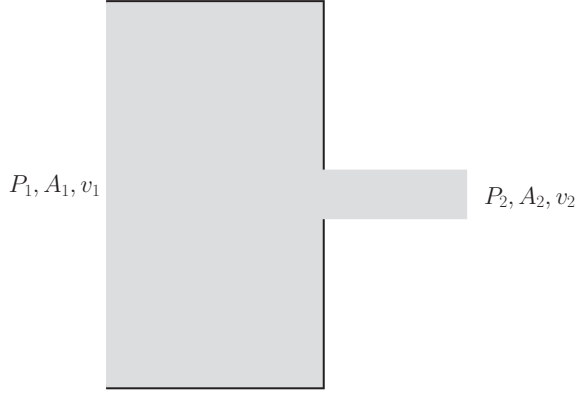


Fig. 8.1: Flow of a fluid through a small orifice. Simple minded application of conservation laws for energy and momentum leads to a paradox. In reality, the cross section of the outgoing stream contracts to avoid this paradox which leads to a phenomenon called Vena Contracta.

the Bernoulli's equation $P + (1/2)\rho v^2 = \text{constant}$ along streamlines, we get the result

$$v_2^2 \approx 2 \frac{P_1 - P_2}{\rho} . \quad (8.23)$$

Let us next try to get the same result using force balance. Since the mass flux across an area is $\rho v A$, the momentum flux is $\rho v^2 A$. The net flux of momentum through the region bound by an area A_1 on the left and area A_2 on the right, is therefore

$$\frac{dp}{dt} = \rho (v_2^2 A_2 - v_1^2 A_1) \approx \rho v_2^2 A_2 , \quad (8.24)$$

when $A_1 v_1^2 \ll A_2 v_2^2$. This rate of change of momentum is caused by the net force on the volume given by

$$F \approx P_1 A_1 - [P_1 (A_1 - A_2) + P_2 A_2] = (P_1 - P_2) A_2 . \quad (8.25)$$

Equating the force in Eq. (8.25) to the rate of change of momentum in Eq. (8.24) we get the result

$$v_2^2 \approx \frac{P_1 - P_2}{\rho} , \quad (8.26)$$

which rudely contradicts the result in Eq. (8.23). Clearly, energy conservation cannot contradict momentum conservation ? Where did we go wrong?

*... conflicts
with that from
momentum
conservation!*

Interestingly enough, the logic and the analysis based on Fig. 8.1 is quite correct but the figure itself is wrong! Nature, which knows that both energy and momentum need to be conserved, adapts to the situation by making the cross section of the outgoing stream contract as it flows. This phenomena called “Vena Contracta” was (probably) first discussed by Torricelli. To see how this works out, assume that the pressure, area and velocity changes from the values (P_2, A_2, v_2) to (P_3, A_3, v_3) as the stream proceeds with $P_3 \ll P_1$. In this case, we get the momentum flux as

*Nature knows
physics*

$$\frac{dp}{dt} = \rho (v_3^2 A_3 - v_1^2 A_1) \approx \rho v_3^2 A_3 \approx 2P_1 A_3, \quad (8.27)$$

where we have used Bernoulli's equation with $P_3 \ll P_1$. The force needed to cause this momentum change is now given by

$$F \approx P_1 A_1 - [P_1 (A_1 - A_2) + P_3 A_3] = (P_1 A_2 - P_3 A_3) \approx P_1 A_2. \quad (8.28)$$

The force balance now leads to the area contraction:

$$A_3 = \frac{A_2}{2}, \quad (8.29)$$

which will save the situation. This is, of course, a rather crude estimate and observations suggest a value close to 0.64 rather than 0.5 which we have obtained. But the basic physics of the problem is indeed what we have described.

One can model the 2-dimensional flow in this case using the fact that the real and complex parts of any analytic function satisfy the Laplace equation. With a clever choice of such functions, one can obtain an analytical model in which the contraction factor is $\pi(2 + \pi)^{-1} \approx 0.61$. Such a modeling also shows that nearly 90 per cent of the contraction occurs within a distance which is about 0.4 of the width of the orifice.

Isochronous Curiosities: Classical and Quantum

9

Your study of classical mechanics usually begins with the analysis of a particle of mass m moving in one dimension under the action of a potential $V(x)$. This is probably the simplest problem in classical mechanics and possibly the whole of physics. As we shall see, this apparent simplicity is rather deceptive and this problem hides some interesting surprises [34].

*The simplest
problem in physics,
or is it?*

Using the constancy of the total energy, $E = (1/2)m\dot{x}^2 + V(x)$, one can write down the equation determining the trajectory of the particle $x(t)$ in the form of the integral

$$t(x) = \sqrt{\frac{m}{2}} \int^x \frac{dx}{\sqrt{E - V(x)}}. \quad (9.1)$$

For a given $V(x)$, this determines the inverse function $t(x)$ and the problem is completely solved. In this chapter, we are interested in the case of bounded oscillations of a particle in a potential well $V(x)$ which has the general shape like the one shown in Fig. 9.1. The potential has a single minimum and increases without bound as $|x| \rightarrow \infty$. For a given value of energy E , the particle will oscillate between the two turning points $x_1(E)$ and $x_2(E)$ which are given by the roots of the equation $V(x) = E$. The period of oscillation between the two turning points can be immediately written down using Eq. (9.1) as:

*The period of
oscillation*

$$T(E) = \sqrt{\frac{m}{2}} \int_{x_1(E)}^{x_2(E)} \frac{dx}{\sqrt{E - V(x)}}. \quad (9.2)$$

(This is actually one-half of the time it takes for the particle to return to the original position; but we will call it period for simplicity.) For a general potential $V(x)$, the result of integration on the right hand side will depend on the value of the energy E . In other words, the period of oscillation will depend on the energy of the particle; equivalently, if one imagines releasing the particle from rest at the location $x = x_1$, say, then the period will depend on the amplitude x_1 of oscillation.

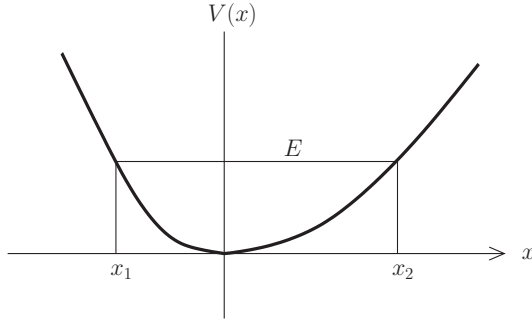


Fig. 9.1: A one-dimensional potential with a single minimum which supports oscillations

For a simple class of potentials, it is quite easy to determine how the period T scales with the energy E . Consider, for example, a class of potentials of the form $V(x) = kx^{2n}$ where n is an integer. These potentials are symmetric in the x -axis and have a minimum at $x = 0$ with the minimum value being $V_{\min} = 0$. In this case, by introducing a variable q such that $q = (k/E)^{1/2n}x$ the energy dependence of the integral in Eq. (9.2) can be easily identified to give

$$T(E) \propto \frac{1}{\sqrt{E}} E^{1/2n} \int_0^1 \frac{dq}{\sqrt{1-q^{2n}}} \propto E^{\frac{1}{2}(\frac{1-n}{n})}. \quad (9.3)$$

Again, harmonic oscillator seems special ...

For all values of n other than $n = 1$, the period T has a non-trivial dependence on the energy. However, when $n = 1$, which corresponds to the harmonic oscillator potential $V(x) = kx^2$, the period is independent of the energy. This, of course, is the well known result that the period of a harmonic oscillator does not depend on the amplitude of the oscillator. The above analysis also shows that amongst all the *symmetric* potentials of the form $V(x) \propto x^{2n}$, *only* the harmonic oscillator has this property.

... but is it, really?

Let us now consider the inverse problem. Suppose you are given the function $T(E)$. Is it possible to determine the potential $V(x)$? For example, if the period is independent of the amplitude, what can we say about the form of the potential $V(x)$? Should it necessarily be a harmonic oscillator potential or can it be more general?

Before launching into a mathematical analysis, let me describe a simple example which deserves to be better known than it is. Consider a potential of the form

$$V(x) = ax^2 + \frac{b}{x^2} \quad (a > 0, b \geq 0), \quad (9.4)$$

in the region $x > 0$. In this region, the potential has a distinct minimum at $x_{\min} = (b/a)^{1/4}$ with the minimum value of the potential being $2\sqrt{ab}$.

The potential is symmetric in x and hence has two minima in the full range $-\infty < x < \infty$; but we shall confine our attention to the range $x > 0$. By shifting the origin suitably we can make the potential in this range to look like the one in Fig. 9.1. For any finite energy, a particle will execute periodic oscillations in this potential. It turns out that *the period of oscillation in this potential is independent of the amplitude* just as in the case of a harmonic oscillator potential! So clearly, a harmonic oscillator is not unique in having this property.

A rival to the oscillator

There are several ways to prove this result. The most difficult one involves evaluating the integral in Eq. (9.2) with $V(x)$ given by Eq. (9.4). The cutest procedure is probably the following. Consider a particle moving, not in one dimension but in two (say in the xy plane), under the action of a two dimensional harmonic oscillator potential

A simple trick

$$V(x, y) = \frac{1}{2}m\omega^2(x^2 + y^2). \quad (9.5)$$

Clearly, under the influence of such a potential, the particle will oscillate with a period which is independent of its energy. Now consider the same problem in polar coordinates instead of Cartesian coordinates. The conservation of energy now gives

$$E = \frac{1}{2}m(\dot{x}^2 + \dot{y}^2) + \frac{1}{2}m\omega^2(x^2 + y^2) = \frac{1}{2}m(\dot{r}^2 + r^2\dot{\theta}^2) + \frac{1}{2}m\omega^2r^2. \quad (9.6)$$

Using the fact that for such a motion — under the central force $V(r) \propto r^2$ — the angular momentum $J = mr^2\dot{\theta}$ is conserved, this expression can be rewritten in the form

$$E = \frac{1}{2}m\dot{r}^2 + \frac{1}{2}m\omega^2r^2 + \frac{1}{2}\frac{J^2}{mr^2} = \frac{1}{2}m\dot{r}^2 + Ar^2 + \frac{B}{r^2}, \quad (9.7)$$

with $A = (1/2)m\omega^2$, $B = J^2/2m$. Mathematically, this is identical to the problem of a particle moving in one dimension under the action of a potential of the form in Eq. (9.4). But we know by construction that the period of oscillation does not depend on the conserved energy E in the case of Eq. (9.7). It follows that the potential in Eq. (9.4) must also have this property.

Since the potential in Eq. (9.4) depends on two parameters a and b you might have thought that the frequency of oscillation ω_0 will also depend on both a and b . The first surprise is that $\omega_0 = 2(2a/m)^{1/2}$ which is independent of b ! This result is most easily found by using the fact that — because the frequency is independent of the amplitude — it must be the same as that for very small oscillations near the minimum. Near the minimum at $x_m = (b/a)^{1/4}$, the potential has the form $V(x) = 2\sqrt{ab} + 4a(x - x_m)^2$ leading to the above result.

*Surprise 1:
Period is
independent
of b !*

Surprise 2: But you can't get it by setting $b = 0$!

Next, once you are told ω_0 is independent of b , you might think it must be the value for the potential obtained by setting $b = 0$, which is $(2a/m)^{1/2}$. The second surprise is that this guess is also not correct! This is because, however small b may be, the potential does rise to infinity near origin and the term (b/x^2) dominates near $x = 0$. (You can think of this as an infinite barrier at $x = 0$ in the limiting case which will double the frequency and halve the period.) So the net effect of b is *only* to double the frequency of oscillation from $(2a/m)^{1/2}$ (when $b = 0$) to $2(2a/m)^{1/2}$ (when $b \neq 0$).

Isochronous potentials, defined

Potentials like that of the harmonic oscillator, or the one in Eq. (9.4) are called isochronous potentials, with the term referring to the property that the period is independent of the amplitude. It is not difficult to see that there are actually an infinite number of such potentials. In fact, for every function $T(E)$, one can construct an infinite number of potentials $V(x)$ such that Eq. (9.4) holds. We will now describe [35] an elementary way to construct them.

New potential from the old

We begin by noting that the period $T(E)$ is determined by the integral in Eq. (9.2), which is essentially the area under the curve $(E - V(x))^{-1/2}$. Consider a potential $V_1(x)$ for which the energy dependence of the period is given by a function $T(E)$. Let us now construct another potential $V_2(x)$ by “shearing” the original potential $V_1(x)$ parallel to x -axis. This is done by shifting the potential curve horizontally by an amount $\Delta(V)$ at every value of V using some arbitrary function $\Delta(V)$. The only restriction on the function $\Delta(V)$ is that the resulting potential should be single valued everywhere. A moment of thought shows that such a shift leaves the area under the curve invariant and hence $T(E)$ does not change. In other words, given any potential $V(x)$, there are infinite number of other potentials for which you will get the same period-energy dependence $T(E)$; each of these potentials are determined by the choice of the ‘shearing’ function $\Delta(V)$.

Worked out example

In the case of a harmonic oscillator potential, the distance $h(V)$ between the two turning points (“width”) varies as \sqrt{V} when the potential is measured from its minima. Since Eq. (9.4) has the isochronous property, we would suspect that it is probably obtained from the harmonic oscillator potential by a shearing motion keeping the width $h(V)$ varying as $(V - V_{\min})^{1/2}$. This is indeed true and we can demonstrate it as follows. From Eq. (9.4), we can determine the inverse, double valued function $x(V)$ through the equation

$$ax^4 + b - Vx^2 = 0. \quad (9.8)$$

If the roots of this equation are x_1^2 and x_2^2 , we immediately have $x_1^2 + x_2^2 = V/a$ and $x_1^2 x_2^2 = b/a$. Elementary algebra now gives

$$h(V)^2 = (x_1 - x_2)^2 = \frac{V}{a} - 2\sqrt{\frac{b}{a}}. \quad (9.9)$$

Or, equivalently,

$$h(V) = \frac{1}{\sqrt{a}} (V - V_{\min})^{1/2} . \quad (9.10)$$

This shows that the potential in Eq. (9.4) is indeed obtained by a shearing of the harmonic oscillator potential.

If you do not like such a geometric argument, here is a more algebraic derivation of the same result [36]. Let us suppose that we are given the function $T(E)$ and are asked to determine the potential $V(x)$ which is assumed to have a single minima and a shape roughly like the one in Fig. 9.1. We can always arrange the coordinates such that the minimum of the potential lies at the origin of the coordinate system. The shape of the curve in the regions $x > 0$ and $x < 0$ will, of course, be different. In order to maintain single valuedness of the inverse function $x(V)$, we will denote the function as $x_1(V)$ in the region $x < 0$ and $x_2(V)$ in the region $x > 0$. Once this is done, we can replace dx in the integral in Eq. (9.2) by $(dx/dV)dV$. This allows us to write

Same result, from algebra

$$T(E) = \sqrt{2m} \int_0^E \left[\frac{dx_2}{dV} - \frac{dx_1}{dV} \right] \frac{dV}{\sqrt{E-V}} \equiv \frac{1}{\pi} \int_0^E \frac{dF}{dV} \frac{dV}{\sqrt{E-V}} , \quad (9.11)$$

where $F(V) \equiv \pi\sqrt{2m}[x_2(V) - x_1(V)]$. This is an integral equation (called Abel's integral equation) which, fortunately, can be inverted by a standard trick. One can easily show that if

Try it out!

$$\frac{1}{\pi} \int_a^t \frac{df}{dx} \frac{dx}{\sqrt{t-x}} = Q(t) , \quad (9.12)$$

then

$$f(x) - f(a) = \int_a^x Q(t) \frac{dt}{\sqrt{x-t}} . \quad (9.13)$$

Using this result and noting that, in our case, $a = 0$ and $F(0) = 0$, we get the final result

$$x_2(V) - x_1(V) = \frac{1}{\pi\sqrt{2m}} \int_0^V \frac{T(E) dE}{\sqrt{V-E}} . \quad (9.14)$$

This result shows explicitly that the function $T(E)$ can determine only the “width” of the curve $x_2(V) - x_1(V)$. The family of curves which has the same width will give rise to the same $T(E)$ and vice-versa. The shearing motion by which we transform one potential to another preserves this width and hence the functional form of $T(E)$.

So far, we have explored the *classical* properties of potentials in which the period of oscillation of a particle is independent of the amplitude. A natural question to ask will be whether these potentials exhibit any inter-

What happens in QM?

esting behaviour in the *quantum mechanical* context. We will now look at some quantum peculiarities [37] of the isochronous potentials.

In quantum theory, the potentials like the one in Fig. 9.1 will have a set of discrete energy levels E_n . Formally inverting the function $E(n)$ — which is originally defined only for integral values of n — one can obtain the inverse function $n(E)$ for this system. This function essentially plays the role analogous to $T(E)$ in the case of quantum theory. We can now ask whether one can determine the potential $V(x)$ given the energy levels E_n or, equivalently, the function $n(E)$. It turns out that one can do this fairly easily in the *semi-classical* limit corresponding to large n . To see this, recall that the energy E_n of the n -th level of a quantum mechanical system is given by the Bohr quantization condition

The semi-classical limit

$$\pi n(E) \simeq \frac{1}{\hbar} \int_{x_1}^{x_2} p dx = \sqrt{\frac{2m}{\hbar^2}} \int_{x_1}^{x_2} \sqrt{E - V} dx. \quad (9.15)$$

(To be precise the n in the left hand side should be $[n - (1/2)]$, but we will work with n ; you can think of this as the $n \gg 1$ limit.) If we differentiate both sides of this equation with respect to E , we get:

$$\sqrt{\frac{2\hbar^2}{m}} \frac{dn}{dE} = \frac{1}{\pi} \int_{E_0}^E \frac{dx}{dV} \frac{dV}{\sqrt{E - V}}, \quad (9.16)$$

where E_0 is the solution to the equation $n'(E_0) = 0$, so that both sides vanish at $E = E_0$. Again using Eq. (9.12) and Eq. (9.13), we get:

$$\begin{aligned} x(V) - x(E_0) &= \sqrt{\frac{2\hbar^2}{m}} \int_{E_0}^V \frac{dn}{dE} \frac{dE}{\sqrt{V - E}} \\ &= \sqrt{\frac{2\hbar^2}{m}} \int_{n(E_0)}^{n(V)} \frac{dn}{\sqrt{V - E(n)}}. \end{aligned} \quad (9.17)$$

The limits of integration are obtained by inverting the function $E(n)$ to get $n(E)$ and substituting the values. This determines the form of the potential $V(x)$ — in terms of the inverse function $x(V)$ — such that in the semi-classical limit it will have the energy levels given by the function $E(n)$.

Though we obtained the above result for a one-dimensional motion with a Cartesian x -axis, it is obvious that a similar formula should be applicable for energy levels in a spherically symmetric potential $V(r)$ provided we only consider the zero angular momentum quantum states. As a curiosity, consider the potential which will reproduce the energy levels that vary as n^{-2} , which — as we know — arises in the case of the Coulomb problem:

$$E_n = -\frac{me^4 Z^2}{2\hbar^2 n^2} \equiv -\frac{C}{n^2}. \quad (9.18)$$

Try it out for the Hydrogen atom

In this case we can take $E_0 = -\infty$ since $n'(-\infty) = 0$. This also gives the lower limit on integration in Eq. (9.17) to be $n(E_0) = n(-\infty) = 0$ and $r(E_0) = r(-\infty) = 0$. An elementary integration of Eq. (9.17) will give

$$\sqrt{\frac{m}{2\hbar^2}} r(V) = \int_0^{n(V)} \frac{ndn}{\sqrt{C + Vn^2}} = \frac{1}{V} (Vn^2 + C)^{1/2} \Big|_0^{n(V)}. \quad (9.19)$$

The contribution from the upper limit vanishes since $n^2(V) = -C/V$ and the lower limit gives $-\sqrt{C}/V$ so that we get the result

$$r = -\frac{Ze^2}{V}; \quad V(r) = -\frac{Ze^2}{r}, \quad (9.20)$$

which, of course, we know is exact. This is one of the many curiosities in the Coulomb problem — viz. the semi-classical result is actually exact — and could be added to the list in Chapter 4. (However, we cheated a little bit in this case; see Box 9.1)

What! Semi-classical result is exact?!

Box 9.1: The Langer trick

The result obtained in Eq. (9.20) suggests that, if we calculate the energy levels in the $(-1/r)$ potential by the WKB approximation, we get the correct result that $E_n \propto -(1/n^2)$. But to do this, we have implicitly set the angular momentum to zero and have looked at the s -states of the atom. If we try to do this properly, we are in for a bit of surprise.

We know that the radial Schrödinger equation for a central potential $V(r)$ corresponding to the angular momentum eigenvalue $\ell(\ell+1)$ is given by

$$\frac{d^2\psi_\ell(r)}{dr^2} + \left(\frac{2m(E - V(r))}{\hbar^2} - \frac{\ell(\ell+1)}{r^2} \right) \psi_\ell(r) = 0; \quad \psi_\ell(0) = 0. \quad (9.21)$$

If we use the standard WKB quantization formula in Eq. (9.15) with n replaced by $(n - 1/2)$ and the WKB momentum being

$$p(r) = \left[2m(E - V(r)) - \ell(\ell+1) \frac{\hbar^2}{r^2} \right]^{1/2} \quad (9.22)$$

and $V(r) = -Ze^2/r$, we find that the energy levels are given by

$$E_{p\ell}^{WKB} = \frac{-mZ^2e^4}{2\hbar^2 \left\{ n - 1/2 + [\ell(\ell+1)]^{1/2} \right\}^2}; \quad n = 1, 2, 3, \dots \quad (9.23)$$

If you do it right, you get it wrong!

This is clearly wrong because it says energy levels depend on ℓ and are degenerate for every value of n ! For $\ell = 0$ you get the correct result for $n \gg 1/2$. The correct result should have only n^2 in the denominator.

Normally, one would have let it go at that saying WKB gives the wrong result, except that Langer found an interesting way of getting around this issue. What Langer did was to replace the WKB momentum in Eq. (9.22) by an effective momentum given by

$$p^{\text{eff}}(r) \equiv \left[2m(E - V(r)) - \left(\ell + \frac{1}{2} \right)^2 \frac{\hbar^2}{r^2} \right]^{1/2}. \quad (9.24)$$

That is, he replaced $\ell(\ell + 1)$ by $[\ell + (1/2)]^2$. This corresponds to adding — out of the blue — a potential $\hbar^2/(8mr^2)$. Incredibly enough, if you use p^{eff} in the WKB formula you get the right result. It turns out that this modification extends the validity of the WKB method [38–40] for a wide class of potentials, regular or singular, attractive or repulsive. There are, however, exceptions to this rule which makes the situation either fascinating or unclear based on your point of view!

A little cheating gets the right result!

There is another interesting feature that arises in the quantum theory related to isochronous potentials. It is well known that when we move from classical to quantum mechanics, the harmonic oscillator potential leads to equidistant energy levels. Curiously enough, all the isochronous potentials have this property *in the semi-classical limit*. This is most easily seen by differentiating Eq. (9.15) with respect to E and using Eq. (9.2) so as to obtain

$$\frac{dn}{dE} = \frac{1}{\pi} \sqrt{\frac{m}{2\hbar^2}} \int_{x_1}^{x_2} \frac{dx}{\sqrt{E - V}} = \frac{T(E)}{\pi\hbar}. \quad (9.25)$$

In other words, the quantum numbers are given by the equivalent formula

$$n(E) \simeq \frac{1}{\pi\hbar} \int T(E) dE, \quad (9.26)$$

which nicely complements the first equation in Eq. (9.15). If the potential is isochronous, then $T(E) = T_0$ is a constant independent of E and the integral immediately gives the linear relation between E and n of the form $E = \alpha n + \beta$ where $\alpha = (\pi\hbar/T_0)$. Clearly, these energy levels are equally spaced just as in the case of harmonic oscillators.

In the case of the potential in Eq. (9.4), something more surprising happens: The *exact* solution to the Schrödinger equation itself has equally spaced energy levels! I will indicate briefly how this analysis proceeds leaving out the algebraic details. To begin with, we can redefine the po-

In the semiclassical limit, all isochronous potentials have equally spaced energy levels

Rivalling the harmonic oscillator, again!

tential to the form

$$V(x) = \left[Ax - \frac{B}{x} \right]^2; \quad A^2 \equiv a, \quad B^2 \equiv b, \quad (9.27)$$

by adding a constant so that the minimum value of the potential is zero at $x = (B/A)^{1/2}$. The frequency of oscillations in this potential is $\omega_0 = (8a/m)^{1/2}$. To study the Schrödinger equation for the potential in Eq. (9.27), it is convenient to introduce the usual dimensionless variables $\xi = (m\omega_0/\hbar)^{1/2}x$, $\varepsilon = 2E/(\hbar\omega_0)$ and $\beta = B(2m)^{1/2}/\hbar$, in terms of which the Schrödinger equation takes the form:

$$\psi'' + \left[\varepsilon - \left(\frac{1}{2}\xi - \frac{\beta}{\xi} \right)^2 \right] \psi = 0. \quad (9.28)$$

As $\xi \rightarrow \infty$, the β/ξ term becomes negligible and — as in the case of standard harmonic oscillator — the wavefunctions will die as $\exp[-(1/4)\xi^2]$. Near the origin, the Schrödinger equation can be approximated as $\xi^2\psi'' \approx \beta^2\psi$ which has solutions of the form $\psi \propto \xi^s$ with s being the positive root of $s(s-1) = \beta^2$. We now follow the standard procedure and write the wavefunction in the form $\psi = \phi(\xi)[\xi^s \exp(-(1/4)\xi^2)]$ and look for power law expansion for ϕ of the form

$$\phi(\xi) = \sum_{n=0}^{\infty} c_n \xi^n. \quad (9.29)$$

Substituting this form into the Schrödinger equation will lead, after some algebra, to the recurrence relation

$$\frac{c_{n+2}}{c_n} = \frac{n+s-\varepsilon-\beta+(1/2)}{(n+2)(n+2s+1)}. \quad (9.30)$$

Asymptotically, this will lead to the behaviour $c_{n+2}/c_n \simeq (1/n)$ so that $\phi(\xi) \simeq \exp[(1/2)\xi^2]$ making ψ diverge unless the series terminates. So, ε must be so chosen that the numerator of Eq. (9.30) vanishes for some value of n . Clearly, only even powers of ξ appear in $\phi(\xi)$ allowing us to write $n = 2k$ where k is an integer. Putting everything back, the energy of the k -th level can be written in the form

$$E_k = (k+C)\hbar\omega_0; \quad C = \frac{1}{2} \left[1 - \beta + \left(\beta^2 + \frac{1}{4} \right)^{1/2} \right], \quad (9.31)$$

showing that the energy levels are equally spaced with the width $\hbar\omega_0$ but with C replacing $(1/2)$ in the case of harmonic oscillator.

Once again there are surprises in store for the limit of $\beta = 0$ when we get $C = 3/4$; shouldn't it be $(1/2)$ in this limit? No. As in the classi-

The last surprise

cal case, we have to imagine an infinite barrier at $x = 0$. If the barrier is removed we get back the normal oscillator but with frequency $(1/2)\omega_0$. (Recall that the isochronous potential leads to twice the frequency of the $b = 0$ case.). The energy levels would have been $(n + (1/2))(1/2)\hbar\omega_0$. But the barrier at $x = 0$ requires the wavefunction to vanish there and hence we can only have odd n eigenfunctions. If we set $n = 2k + 1$ the energy levels become $[2k + (3/2)](1/2)\hbar\omega_0 = [k + (3/4)]\hbar\omega_0$ which is the origin of $C = 3/4$!

Do all isochronous potentials lead to evenly spaced energy levels as exact solutions to Schrödinger equation rather than only in the asymptotic limit? The answer is “no”. The simple counter-example is provided by two parabolic wells connected together smoothly at the minima with $V(x) = (1/2)m\omega_R^2x^2$ for $x \geq 0$ and $V(x) = (1/2)m\omega_L^2x^2$ for $x \leq 0$. It is obvious that this potential is isochronous classically. Solving the Schrödinger equation requires some effort because you need to ensure continuity of ψ and ψ' at the origin. This leads to a set of energy levels which need to be solved numerically. One then finds that the energy levels are not equally spaced but the departure from even spacing is surprisingly small. There is no simple characterization of potentials which lead to evenly spaced energy levels in quantum theory.

*A cute conjecture
is killed by a cruel
counterexample*

Most courses in electrostatics begin by studying the Gauss law and its application to determine the electric fields produced by simple charge distributions. In this chapter, we revisit one of these problems, viz., the field produced by an infinitely long, straight, line of charge with a constant charge density. As usual, we will do it in a slightly different manner compared to the text books and get ourselves all tied up in knots [41].

Consider an infinite straight line of charge located along the y -axis with a charge density per unit length being λ . We are interested in determining the electric field everywhere due to this line charge. The standard solution to this problem is very simple. We first argue, based on the symmetry, that the electric field at any given point is in the $x-z$ plane and depends only on the distance from the line charge. So we can arrange the coordinate system such that the point at which we want to calculate the field is at $(x, 0, 0)$. If we now enclose the line charge by an imaginary concentric cylindrical surface of radius x and length L , the outward flux of the electric field through the surface is $2\pi xLE$ which should be equal to (4π) times the charge enclosed by the cylinder which is $(4\pi L\lambda)$. This immediately gives $E = (2\lambda/x)$. Dimensionally, the electric field is the charge divided by the square of the length, and since λ is charge per unit length, everything is fine.

The standard result, if you compute \mathbf{E} directly

We will now do it differently and in — what should be — an equivalent way. We compute the electrostatic potential ϕ at $(x, 0, 0)$ due to the line charge along the y -axis and obtain the electric field by differentiating ϕ . Obviously, the potential $\phi(x)$ can only depend on x and λ and must have the dimension of charge per unit length. If we take $\phi \sim \lambda^n x^m$, dimensional analysis immediately gives $n = 1$ and $m = 0$ so that $\phi(x) \propto \lambda$ and is independent of x ! The potential is a constant and the electric field vanishes! We are in trouble.

Serious trouble, if you compute ϕ

Computation of the potential from first principles makes matters *worse*! An infinitesimal amount of charge $dq = \lambda dy$ located between y and $y + dy$ will lead to an electrostatic potential dq/r at the field point where

$r = (x^2 + y^2)^{1/2}$. So the total potential is given by

$$\phi(x) = \lambda \int_{-\infty}^{+\infty} \frac{dy}{\sqrt{x^2 + y^2}} = 2\lambda \int_0^{+\infty} \frac{dy}{\sqrt{x^2 + y^2}}. \quad (10.1)$$

Changing variables from y to $u = y/x$, the integral becomes

$$\phi(x) = 2\lambda \int_0^{+\infty} \frac{du}{\sqrt{1 + u^2}}. \quad (10.2)$$

This result is clearly independent of x and hence a constant (which is what dimensional analysis told us). Much worse, it is an *infinite* constant since the integral diverges at the upper limit. What is going on in such a simple, classic, textbook problem?

To get a sensible result, let us try cutting off the integral in Eq. (10.1) at some length scale $y = \Lambda$. (You may think of the infinite line charge as the limit of a line charge of length 2Λ with $\Lambda \gg x$.) Using the substitution $y = x \sinh \theta$ and taking the limit $\Lambda \gg x$ we get

$$\phi(x) = 2\lambda \int_0^\Lambda \frac{dy}{\sqrt{x^2 + y^2}} = 2\lambda \sinh^{-1} \left(\frac{\Lambda}{x} \right) \approx -2\lambda \ln \left(\frac{x}{2\Lambda} \right), \quad (10.3)$$

where we have used $\Lambda \gg x$ in arriving at the final equality. This potential does diverge when $\Lambda \rightarrow \infty$. But note that the physically observable quantity, the electric field $\mathbf{E} = -\nabla\phi$ is independent of the cut-off parameter Λ and is correctly given by $E_x = 2\lambda/x$. By introducing a cut-off, we seem to have saved the situation.

We can now clearly see what is going on. As the title of this chapter implies, the problem has to do with logarithms which allows a dimensionless function like $\ln(x/2\Lambda)$ to slip into the electrostatic potential without the electric field depending on the arbitrary scale Λ . This requires additivity on the Λ dependence; that is, we need a function $f(x/\Lambda)$ which will reduce to $f(x) + f(\Lambda)$. Clearly, only a logarithm will do.

Once we know what is happening we can figure out other ways of getting a sensible answer. One can, for example, obtain this result from a more straightforward scaling argument by concentrating on the potential difference $\phi(x) - \phi(a)$ where a is some arbitrary scaling distance we introduce into the problem. From dimensional analysis, it follows that the potential difference must have the form $\phi(x) - \phi(a) = \lambda F(x/a)$ where F is a dimensionless function. Evaluating this expression for $a = 1$, say, in some units, we get $\lambda F(x) = \phi(x) - \phi(1)$. Substituting back, we have the relation $\phi(x) - \phi(a) = \phi(x/a) - \phi(1)$. This functional equation has the unique solutions $\phi(x) = A \ln x + \phi(1)$. Dimensional analysis again tells you that $A \propto \lambda$ but, of course, scaling arguments cannot determine the proportionality constant. However, one can compute the potential differ-

*Infinites?
In electrostatics?!*

*A way-out, with
deeper meaning*

ence by the explicit integral

$$\phi(x) - \phi(a) = 2\lambda \int_0^\infty dy \left(\frac{1}{\sqrt{x^2 + y^2}} - \frac{1}{\sqrt{a^2 + y^2}} \right). \quad (10.4)$$

We can easily see that this integral is finite. A fairly straightforward calculation leads to:

$$\phi(x) - \phi(a) = -2\lambda \ln(x/a). \quad (10.5)$$

*With hindsight,
you can do it in
many ways*

The numerical value of $\phi(x)$ in this expression is independent of the length scale a introduced in the problem. In that sense the scale of ϕ is determined only by λ which, as we said before, has the correct dimensions. But to ensure finite values for the expressions, we need to introduce an arbitrary length scale a which is the key feature I want to emphasize in this discussion.

It turns out that such phenomena, in which naive scaling arguments break down due to the occurrence of the logarithmic function, is a very general feature in several areas of physics especially in the study of the renormalization group in high energy physics. What we have here is a very elementary analogue of this result. In all these cases, we introduce a length scale into the problem to make some unobservable quantities (like the potential) finite but arrange matters such that *observable* quantities remain independent of this scale which we throw in.

If you thought this was too simple, here is a more sophisticated occurrence of a logarithm for similar reasons.

Consider the Schrödinger equation in *two* dimensions for an attractive Dirac delta function potential $V(\mathbf{x}) = -V_0 \delta(\mathbf{x})$ with $V_0 > 0$. The vector \mathbf{x} is in two dimensional space, and we look for a stationary bound state wavefunction $\psi(\mathbf{x})$ which satisfies the equation

*Schrödinger
equation for
delta function
potential*

$$\left(-\frac{\hbar^2}{2m} \nabla^2 - V_0 \delta(\mathbf{x}) \right) \psi(\mathbf{x}) = -|E| \psi(\mathbf{x}), \quad (10.6)$$

where $-|E|$ is the negative bound state energy. We rescale the variables by introducing $\lambda = 2mV_0/\hbar^2$ and $\mathcal{E} = 2m|E|/\hbar^2$, so that this equation reduces to

$$(\nabla^2 + \lambda \delta(\mathbf{x})) \psi(\mathbf{x}) = \mathcal{E} \psi(\mathbf{x}). \quad (10.7)$$

Everything up to this point could have been done in any spatial dimension. In D -dimensions, the Dirac delta function $\delta(\mathbf{x})$ has the dimension L^{-D} . The kinetic energy operator ∇^2 , on the other hand, always has the dimension L^{-2} . This leads to a peculiar behaviour when $D = 2$. We find that, in this case, λ is dimensionless while \mathcal{E} has the dimension of L^{-2} . Since the scaled binding energy \mathcal{E} has to be determined entirely in terms of the parameter λ , we have a serious problem in our hands. There is no

*Trouble in store
for $D = 2$*

way we can determine the form of \mathcal{E} without a dimensional constant — which we do not have!

We now solve Eq. (10.7) to see the manifestation of this problem more clearly. This is fairly easy to do by Fourier transforming both sides and introducing the momentum space wavefunction $\phi(\mathbf{k})$ by

$$\phi(\mathbf{k}) = \int d^2\mathbf{x} \psi(\mathbf{x}) \exp(-i\mathbf{k} \cdot \mathbf{x}) . \quad (10.8)$$

The left hand side of Eq. (10.7) leads to the term $[-k^2\phi(\mathbf{k}) + \lambda\psi(0)]$ while the right hand side gives $\mathcal{E}\phi(\mathbf{k})$. Equating the two, we get:

$$\phi(\mathbf{k}) = \frac{\lambda\psi(0)}{k^2 + \mathcal{E}} . \quad (10.9)$$

We now integrate this equation over all \mathbf{k} . The left hand side will then give $(2\pi)^2\psi(0)$ which can be canceled out on both sides because $\psi(0) \neq 0$. (This is, of course, needed for $\phi(\mathbf{k})$ in Eq. (10.9) to be non-zero.) We then get the result

$$\frac{1}{\lambda} = \frac{1}{4\pi^2} \int \frac{d^2\mathbf{k}}{k^2 + \mathcal{E}} = \frac{1}{4\pi^2} \int \frac{d^2\mathbf{s}}{s^2 + 1} . \quad (10.10)$$

*Disaster, derived
explicitly*

The second equality is obtained by changing the integration variable to $\mathbf{s} = \mathbf{k}/\sqrt{\mathcal{E}}$. This equation is supposed to determine the binding energy \mathcal{E} in terms of the parameter in the problem λ but the last expression shows that the right hand side is independent of \mathcal{E} ! This is similar to the situation in the electrostatic problem in which we got the integral in Eq. (10.2) which was independent of x . In fact, just as in the electrostatic case, the integral on the right hand side diverges, confirming our suspicion. Of course, we already know that determining \mathcal{E} in terms of λ is impossible due to dimensional mismatch.

One can, at this stage, take the point of view that the problem is simply ill-defined and one would be quite correct. The Dirac delta function, in spite of the nomenclature, is strictly not a function but is, what mathematicians will call, a distribution. It is defined as a limit of a sequence of functions. For example, suppose we consider a sequence of potentials

The easy way out

$$V(\mathbf{x}) = -V_0 \left[\frac{1}{2\pi\sigma^2} \exp(-|\mathbf{x}|^2/2\sigma^2) \right] , \quad (10.11)$$

where \mathbf{x} is a 2D vector and σ is a parameter with the dimension of length. In this case, we will again get Eq. (10.7) but with the Dirac delta function replaced by the Gaussian in the square brackets in Eq. (10.11). But now we have a parameter σ with the dimensions of length and one can imagine the binding energy being constructed out of this. When we take the limit $\sigma \rightarrow 0$, the potential in Eq. (10.11) reduces to a Dirac delta function. This is

what is meant by saying that the delta function is defined as a limiting case of sequence of functions. Here, the functions are Gaussians in Eq. (10.11) parametrized by σ . When we take the limit of $\sigma \rightarrow 0$ the function reduces to the delta function. The trouble is that, when we let σ go to zero, we lose the length scale in the problem and we do not know how to fix the binding energy. Of course, there is no assurance that if one solves a differential equation with an input function $V(\mathbf{x}; \sigma)$ which depends on a parameter σ and take a somewhat dubious limit of $\sigma \rightarrow 0$, the solutions will have a sensible limit. So one can say that the problem is ill-defined.

Rather than leaving it at that, we can attempt something similar to what we did in the electrostatic case. Evaluating the integral in Eq. (10.10) with a cut-off at some value $k_{max} = \Lambda$ with $\Lambda^2 \gg \mathcal{E}$, we get

Let us be adventurous

$$\frac{1}{\lambda} = -\frac{1}{4\pi} \ln \left(\frac{\mathcal{E}}{\Lambda^2} \right), \quad (10.12)$$

which can be inverted to give the binding energy to be:

$$\mathcal{E} = \Lambda^2 \exp(-4\pi/\lambda), \quad (10.13)$$

where the scale is fixed by the cut-off parameter. Of course this is similar to what we would have got if we actually used a potential with a length scale.

One way of interpreting this result is by taking a clue from what is done in quantum field theory. The essential idea is to accept up front that the theory requires an extra scale with proper dimensions for its interpretation. We then treat the coupling constant as a function of the scale at which we probe the system. Having done that, we arrange matters so that the observed results are actually independent of the scale we have introduced. In this case, we will define a physical coupling constant by

QFT based insight for QM

$$\lambda_{phy}^{-1}(\mu) = \lambda^{-1} - \frac{1}{4\pi} \ln(\Lambda^2/\mu^2) = -\frac{1}{4\pi} \ln \left(\frac{\mathcal{E}}{\mu^2} \right), \quad (10.14)$$

where μ is an arbitrary but finite scale. Obviously $\lambda_{phy}(\mu)$ is independent of the cut-off parameter Λ . The binding energy is now given by

$$\mathcal{E} = \mu^2 \exp(-4\pi/\lambda_{phy}(\mu)), \quad (10.15)$$

which, in spite of appearance, is independent of the scale μ . This is similar to our Eq. (10.5) in the electrostatic problem, in which we introduced a scale a but $\phi(x)$ was independent of a .

In quantum field theory, the above result will be interpreted as follows: Suppose one performs an experiment to measure some observable quantity (like the binding energy) of the system as well as some of the parameters describing the system (like the coupling constant). If the experiment is performed at an energy scale corresponding to μ (which, for

The running coupling constant: key concept in QFT

example could be the energy of the particles in a scattering cross-section measurement, say), then one will find that the measured value of the coupling constant depends on μ . But when one varies μ in an expression like Eq. (10.15), the variation of λ_{phys} will be such that one gets the same value for \mathcal{E} .

When you think about it, it does make lot of sense. After all, the parameters we use in our equations (like λ_{phy}) as well as some of the results we obtain (like the binding energy \mathcal{E}) need to be determined by suitable experiments. In the quantum mechanical problems one can think of scattering of a particle with momentum k (represented by an incident plane wave, say) by a potential. The resulting scattering cross-section will contain information about the potential, especially the coupling constant λ . If the scattering experiment introduces a (momentum or length) scale μ , then one can indeed imagine the measured coupling constant to be dependent on that scale μ . But we would expect physical predictions of the theory (like \mathcal{E}) to be independent of μ . This is precisely what happens in quantum field theory and the toy model above is a simple illustration.

*Scattering in 2D
delta function
potential*

It is fairly straightforward to see how all these comes about in the case of *scattering* in the 2-dimensional Dirac delta function potential. The analysis is very similar to what was done above and the formalism uses the scattering theory developed in Chapter 4. Let me briefly indicate the key results.

When we study scattering solutions to the Schrödinger equation, we take $E \equiv k^2/2m > 0$ in contrast to the bound state problems in which we assume $E = -|E| < 0$. If you now carry through the analysis similar to the one done above, you will easily find that a scattering state wavefunction will be given by

$$\psi_{\mathbf{k}}(\mathbf{x}) = e^{i\mathbf{k}\cdot\mathbf{x}} + \lambda \psi_{\mathbf{k}}(\mathbf{0})G(\mathbf{x}) , \quad (10.16)$$

where $G(\mathbf{x})$ is the 2-dimensional Green's function given by

$$G(\mathbf{x}) = \int \frac{d^2\mathbf{p}}{(2\pi)^2} \frac{e^{i\mathbf{p}\cdot\mathbf{x}}}{\mathbf{p}^2 - \mathbf{k}^2 - i\epsilon} , \quad (10.17)$$

which can be expressed in terms of the zero-th order Hankel function. Evaluating Eq. (10.16) at the origin now leads to the consistency condition

$$\psi_{\mathbf{k}}(\mathbf{0}) = \frac{1}{1 - \lambda G(\mathbf{0})} , \quad (10.18)$$

which again lands us in trouble because $G(\mathbf{0})$ is logarithmically divergent. As in the previous case, let us evaluate the integral in Eq. (10.17) with a cut-off at $|\mathbf{p}| \leq \mu$. This will lead to

$$G(\mathbf{0}) = \frac{1}{4\pi} \ln \left(\frac{\mu^2}{-k^2} \right) . \quad (10.19)$$

Using Eq. (10.18) and Eq. (10.19) in Eq. (10.16) and using the asymptotic expansion of the Hankel function,

$$H_0^1(kr) \rightarrow \left(\frac{2}{i\pi kr} \right)^{1/2} e^{ikr} \quad (kr \rightarrow \infty), \quad (10.20)$$

you can easily determine the scattering amplitude $f(\theta)$ to be

$$f(\theta) = \sqrt{\frac{2}{\pi k}} \left[\frac{1}{\lambda} - \frac{1}{4\pi} \ln \left(\frac{\mu^2}{k^2} \right) - \frac{i}{4} \right]^{-1}. \quad (10.21)$$

We now see that, if we analytically continue to negative energies by the replacement $k \rightarrow ik$, then $f(\theta)$ possess a pole at

$$k_{\text{phy}}^2 = \mu^2 \exp \left(-\frac{4\pi}{\lambda} \right) = \mathcal{E}, \quad (10.22)$$

which is precisely the bound state energy we obtained in Eq. (10.15). This agrees with the general result in quantum mechanics that the poles of the scattering amplitude at imaginary values of k occur at the bound state energies. More importantly, the scattering cross section is now given by

$$\frac{d\sigma}{d\theta} = |f(\theta)|^2 = \frac{2}{\pi k} \left[1 + \frac{1}{\pi^2} \ln^2 \left(\frac{k^2}{\mathcal{E}} \right) \right]^{-1}, \quad (10.23)$$

which depends on the regularized bound state energy \mathcal{E} . Suppose we determine the scattering cross section at the value of k given by $k = \mu$. This will allow us to determine the physical coupling constant $\lambda_{\text{phy}}(k)$ using Eq. (10.15). This coupling constant will “run” with the energy scale k at which the scattering experiment is performed but this dependence will be such that the bound state energy \mathcal{E} remains constant.

From Eq. (10.7) it can be seen that, in $D = 1$, the coupling constant λ has the dimensions of L^{-1} so there is no difficulty in obtaining $E \propto \lambda^2$. The one dimensional integral corresponding to Eq. (10.10) is convergent and you can easily work this out to fix the proportionality constant to be $1/4$. The logarithmic divergence occurs in $D=2$ which is known as the critical dimension for this problem. The breaking down of naive scaling arguments and the appearance of logarithms are rather ubiquitous in such a case. (There are other fascinating issues in $D \geq 3$ and in the scattering by potentials but that is another story.) The examples discussed here are all explored extensively in the literature and a good starting point will be Refs. [42–47].

The simplest problem in gravity deals with the description of the gravitational field produced by a spherically symmetric distribution of matter around it. In Newtonian gravity, we will describe it using a gravitational potential which falls as $(-1/r)$ everywhere outside the body. In Einstein's theory, one describes gravity as due to the curvature of spacetime. So, to understand, say, the effects of general relativity in the solar system, we need to determine the spacetime geometry around the Sun. The rigorous way of doing this is to solve Einstein's field equations in this specific context. How would you like to get this key result of general relativity rather cheaply?

In this chapter, we will discuss how this important result of general relativity, viz. the description of the gravitational field around a spherical massive body, can be obtained using just the concepts of special relativity [48, 49]. This curious fact allows you to explore a host of physical phenomena including some aspects of black hole physics. The derivation works only for a special class of spherically symmetric models — for reasons which are not completely clear — but considering how easy it is, it deserves to be known much better.

This stunt is performed by experts; do not try this on your own at home !

We start with the fact that general relativity describes gravity as due to the curvature of spacetime. The difference between a flat space and a curved space is encoded in the generalization of Pythagoras theorem for infinitesimally separated points. For example, a flat 2-dimensional surface (say, a plain sheet of paper) allows us to introduce the standard Cartesian coordinates (x, y) such that the distance between infinitesimally separated points can be expressed in the form $dl^2 = dx^2 + dy^2$ which, of course, is just the standard Pythagoras theorem. In contrast, consider the two dimensional surface on a sphere of radius r on which we have introduced two angular coordinates (θ, ϕ) . The corresponding formula will now read $dl^2 = r^2 d\theta^2 + r^2 \sin^2 \theta d\phi^2$. It is not possible to introduce any other set of coordinates on the surface of sphere such that this expression — usually

See Box 11.1 if you want to know why

called the line interval — reduces to the Pythagorean form. This is the difference between a curved space and a flat space.

Move on from space to spacetime and from points to events. In the flat spacetime, in which we use in special relativity, the “Pythagoras theorem” generalizes to the form

$$ds^2 = -c^2 dt^2 + dx^2 + dy^2 + dz^2 . \quad (11.1)$$

The spatial coordinates appear in the standard form and the inclusion of time introduces the all important minus sign. But one can live with it and treat it as a generalization of the formula $dl^2 = dx^2 + dy^2$ to 4-dimensions (with an extra minus sign). But in a curved spacetime, this expression will not hold and the coordinate differentials like $c^2 dt^2, dx^2$ etc. in the interval will get multiplied by functions of space and time. This is just like we using $\sin^2 \theta d\phi^2$ rather than just $d\phi^2$ to describe the curved 2-dimensional surface of a sphere. The precise manner in which such a modification occurs is determined by Einstein’s equation and depends on the distribution of matter in spacetime. Our aim is to find the spacetime around a massive body by using a trick.

Box 11.1: Why is gravity just geometry?

There are three ingredients which lead to the fascinating result that the effects of gravity must be represented in terms of curved spacetime geometry. The first is the principle of equivalence which tells you that in a sufficiently small region of spacetime you cannot distinguish the effects of gravity from those produced by a suitable accelerated frame (see text). The second is the well known result in special relativity, namely, moving clocks run slower compared to stationary ones. Combining these two, one can show that gravitational potential affects the rate of clocks. The final ingredient is the requirement that our description should be valid in any arbitrary coordinate system because one can no longer distinguish effects of gravity from those of accelerated frame, locally. Let me fill in the details of this argument.

Consider a disc rotating with angular velocity Ω about an axis running through the center perpendicular to the disc. Keep one clock at the center of the disc (which does not move) and another at a radius r which moves with a constant *speed* $v = \Omega r$. Special relativity tells you that when the clock at the origin shows a lapse of time $\Delta t(0)$, the clock at radius r will show a lapse of time $\Delta t(r)$ given by (see Eq. (2.40)):

$$\Delta t(r) = \Delta t(0) \left(1 - \frac{v^2}{c^2}\right)^{1/2} = \Delta t(0) \left(1 - \frac{\Omega^2 r^2}{c^2}\right)^{1/2} . \quad (11.2)$$

The three key ingredients

Someone sitting with the clock at r in a closed cabin will feel the centrifugal acceleration, $\Omega^2 r$ which she cannot distinguish from a gravitational acceleration arising from a gravitational potential ϕ such that $-\partial\phi/\partial r = \Omega^2 r$; this leads to a gravitational potential $\phi = -(1/2)\Omega^2 r^2$. Using principle of equivalence we can now re-express Eq. (11.2) as

$$\Delta t(\phi) = \Delta t(0) \left(1 + \frac{2\phi}{c^2}\right)^{1/2}. \quad (11.3)$$

This result tells you that the flow of time depends on the gravitational potential at which the clock is located. If this does not hold, either principle of equivalence or special relativity should fail!

We next note that the line interval in special relativity is of the form $ds^2 = -c^2 dt^2 + d\mathbf{x}^2$. Any clock at rest anywhere in space has the worldline $d\mathbf{x} = 0$ and all such clocks will measure the proper time $d\tau \equiv ds/c = dt$. This, of course, contradicts the result in Eq. (11.3) and hence we need to modify the line interval of special relativity in the presence of a gravitational field. The simplest modification which will take care of the effect of gravity on the clock will be to change ds^2 to the form

$$ds^2 = -c^2 d\tau^2 = -\left(1 + \frac{2\phi}{c^2}\right) c^2 dt^2 + d\mathbf{x}^2. \quad (11.4)$$

Stationary clocks with $d\mathbf{x} = 0$ will now show a time lapse in accordance with Eq. (11.3) which will depend on the potential they are located at.

There is a beautiful way of verifying whether we are on the right track. We know that the action for a particle in special relativity is given by (see Chapter 2)

$$A = -mc^2 \int d\tau. \quad (11.5)$$

Principle of equivalence tells you that in any local region you can go to the freely falling frame in which the special relativity should hold. Therefore, in the presence of a weak gravitational field the action for a particle must have the same form as Eq. (11.5) with $d\tau$ given by the expression in Eq. (11.4). If you work it out, to the lowest order in $(1/c^2)$ — which is necessary because everything we did is only valid for weak gravitational fields described by a Newtonian potential — we find that

$$\begin{aligned} A &= -mc^2 \int d\tau = -mc^2 \int dt \left[1 + \frac{2\phi}{c^2} - \frac{v^2}{c^2}\right]^{1/2} \\ &\cong -mc^2 \int dt \left[1 - \frac{1}{2} \frac{v^2}{c^2} + \frac{\phi}{c^2}\right]. \end{aligned} \quad (11.6)$$

Gravity affects flow of time

From flat to curved

A cross-check

But this is equivalent to using the Lagrangian

$$L = -mc^2 + \frac{1}{2}mv^2 - m\phi, \quad (11.7)$$

which, except for the constant $(-mc^2)$, is precisely what you would have written down for a particle in a Newtonian gravitational potential!

Why is the Lagrangian $K - V$ in a gravitational field?

So you see that the Lagrangian for a particle in a gravitational field has the strange form, of kinetic energy minus potential energy, because gravity affects the rate of flow of clocks; there is no other good reason for this strange combination. (The innocuous looking constant $(-mc^2)$ in Eq. (11.7) has interesting consequences which we will explore in Chapter 15.) We conclude that, in the presence of a weak gravitational field, the form of ds^2 must get modified, at least as regards the g_{00} component.

But we know that the operational definition of spatial distances will use constancy of speed of light and the measurement of clock rates. This means even the spatial length interval will get affected by the gravitational field requiring the modification of *all* the components of g_{ab} from the special relativistic form of η_{ab} . (This effect is not captured in the above analysis because we were working at the lowest order in $(1/c^2)$; then $c^2 dt^2$ dominates over $d\mathbf{x}^2$.) It turns out that, for a proper description of gravity, you actually need to go beyond the description by a single gravitational potential and use all the ten components of g_{ab} . That will lead you to a curved spacetime.

For this, we will begin with a simple idea which you probably know. Gravity obeys the *principle of equivalence*. Consider, for example, a small box ('Einstein's elevator') which is moving in intergalactic space, away from all material bodies, in some direction ("up") with a uniform acceleration g . We will assume that it is propelled by the rocket motors attached to its bottom. Let us compare the results of any physical experiment performed inside such an elevator with those performed inside a similar box which is at rest on Earth's surface where the gravitational acceleration is g . The principle of equivalence tells you that the results of all physical experiments will be the same in these two cases; that is, you cannot distinguish gravity from an accelerated frame within any sufficiently small region.

Not the recommended procedure to verify the principle of equivalence

An immediate consequence of this principle is that you can make gravity go away, within any small region of space, by choosing a suitable frame of reference usually called the freely falling frame. For example, if you jump off from the twentieth floor of a building you will feel completely weightless ('zero gravity') until you crash to the ground. In such a

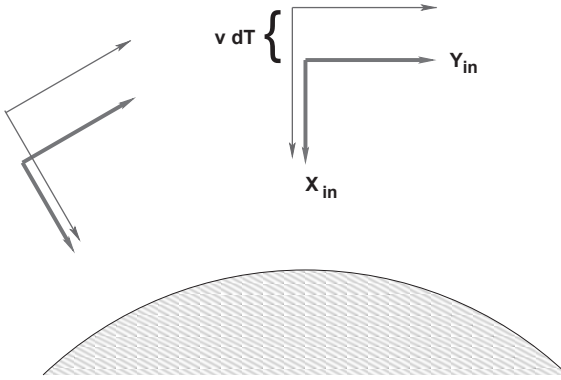


Fig. 11.1: The relation between the freely falling and static coordinate frames around a spherically symmetric body. The thick lines indicate (X_{in}, Y_{in}) axes of the freely falling inertial frame. The thin lines denote the corresponding axes of a static coordinate system glued at a fixed point. (Of course, the figure is not to scale and the coordinate systems are supposed to be infinitesimal in extent!). The radial displacement between the two frames is by the amount $v_r(r)dT$ during an infinitesimal time interval dT . Through every event there is a different freely falling frame related in a definite way to the fixed static frame. The freely falling and static frames coincide at very large distance from the body.

freely falling frame, one can use the laws of special relativity without any problem since gravity is absent.

We want to use this idea to describe the gravitational field of a spherically symmetric body located about the origin. Consider a body of radius R and let us study its gravitational field in the empty space around it (at $r > R$). Let P be a point at a distance r from the origin. If we consider a small box around P which is freely falling towards the origin, then the metric in the coordinates used by a freely falling observer in the box will be just that of special relativity:

$$ds^2 = -c^2 dt_{in}^2 + d\mathbf{r}_{in}^2. \quad (11.8)$$

This is because, in the freely falling frame, the observer is weightless and there is no effective gravity. (The subscript ‘in’ tells you that these are inertial coordinates.) Let us now transform the coordinates from the inertial frame to a frame (T, \mathbf{r}) which will be used by observers who are at rest around the point P . Suppose the freely falling frame is moving with a radial velocity $\mathbf{v}(r)$ around P . To determine this velocity, we can imagine that the freely falling frame started from very large distance from the body with zero velocity at infinity. Then, a simple Newtonian analysis shows that its velocity at P will be $\mathbf{v}(r) = -\hat{\mathbf{r}}\sqrt{2GM/r}$. We now transform from the freely falling inertial frame to the static frame of reference which is glued to the point P using the non-relativistic transformations $dt_{in} = dT$,

Introduce freely-falling-frames around a massive body

Pull a fast one ...

$d\mathbf{r}_{\text{in}} = d\mathbf{r} - \mathbf{v}dT$ between two frames which move with respect to each other with a relative velocity \mathbf{v} . Of course, you have to use infinitesimal quantities in this transformation because you need different freely falling inertial frames at different points, in a non-uniform gravitational field. Also note that $d\mathbf{r}_{\text{in}}$ is *not* an exact differential; you cannot integrate $d\mathbf{r}_{\text{in}} = d\mathbf{r} - \mathbf{v}dT$ to get a coordinate \mathbf{r}_{in} .

What we require is the form of the line element in Eq. (11.8) in terms of the static coordinates. Substituting the transformations in Eq. (11.8), we find the metric in the new coordinates to be

$$ds^2 = - \left[1 - \frac{2GM}{c^2 r} \right] c^2 dT^2 + 2\sqrt{(2GM/r)} drdT + d\mathbf{r}^2. \quad (11.9)$$

... to get the right result!

Incredibly enough, this turns out to be the correct metric describing the spacetime around a spherically symmetric mass distribution of total mass M !

Little bit of cosmetics

As it stands, this line element is not in “diagonal” form in the sense that it has a non-zero $drdT$ term. It will be nicer to have the metric in diagonal form. This can be done by making a coordinate transformation of the time coordinate (from T to t) in order to eliminate the off-diagonal term. We look for a transformation of the form $T = t + Q(r)$ with some function $Q(r)$. This is equivalent to taking $dT = dt + K(r)dr$ with $K = dQ/dr$. Substituting for dT in Eq. (11.9) we find that the off-diagonal term is eliminated if we choose $K(r) = \sqrt{2GM/c^4 r} (1 - 2GM/c^2 r)^{-1}$. In this case, the new time coordinate is:

$$t = \int dT + \frac{1}{c^2} \int dr \frac{\sqrt{(2GM/r)}}{(1 - \frac{2GM}{c^2 r})}. \quad (11.10)$$

The integral in the second term is elementary and working it out you will find that

$$ct = cT - \left[\sqrt{\frac{8GM}{c^2}} r - \frac{4GM}{c^2} \tanh^{-1} \sqrt{\frac{2GM}{c^2 r}} \right]. \quad (11.11)$$

What is more important for us is the final form of the line interval in Eq. (11.9) expressed in the static coordinates with the new time coordinate t . This is given by

$$ds^2 = - \left(1 - \frac{2GM}{c^2 r} \right) c^2 dt^2 + \left(1 - \frac{2GM}{c^2 r} \right)^{-1} dr^2 + r^2 (d\theta^2 + \sin^2 \theta d\phi^2). \quad (11.12)$$

The final result

Some of you might be familiar with this metric, called the *Schwarzschild metric*, which is used extensively both in the study of general relativistic corrections to motion in the solar system and in black hole physics.

Once we have the form of this metric, one can do several things with it using just special relativistic concepts. One simple, but very significant, result which you can obtain immediately is the following. Let us consider a clock which is sitting quietly at some fixed location in space so that, along the clock's worldline, $dr = d\theta = d\phi = 0$. Substituting these into Eq. (11.12) we get the proper time shown by such a clock to be

*Gravitational
redshift*

$$d\tau = \left[1 - \frac{2GM}{c^2 r}\right]^{1/2} dt \equiv \sqrt{|g_{00}(r)|} dt, \quad (11.13)$$

when the coordinate clock time changes by an amount dt . It is obvious from this relation that $d\tau \rightarrow dt$ when $r \rightarrow \infty$. That is, one can think of t as the proper time measured by a clock located far away from the gravitating body.

Consider now an electromagnetic wave train having N crests and troughs which is traveling radially outward from some point \mathbf{x} to an infinite distance away from the central body. An observer located near \mathbf{x} can measure the frequency of the wave train by measuring the time $\Delta\tau$ it takes for the N troughs to cross her and using the result $\omega = N/\Delta\tau$. An observer at large distances will do the same using her clocks. Since the frequency of radiation $\omega(\mathbf{x})$ measured by local observers, as the radiation propagates from event to event in a curved spacetime, is inversely related to the time measured by the local clock, it follows that $\omega(\mathbf{x}) \propto [|g_{00}(\mathbf{x})|]^{-1/2}$. If $g_{00} \approx -1$ at very large distances from a mass distribution, then the frequency of radiation measured by an observer at infinity (ω_∞) will be related to the frequency of radiation emitted at some point \mathbf{x} by

$$\omega_\infty = \omega(\mathbf{x}) \sqrt{|g_{00}(\mathbf{x})|}. \quad (11.14)$$

Another use of this metric is to study orbits of particles around a massive body. The formal way of doing this is to use the Hamilton-Jacobi equation, $g^{ab}\partial_a S \partial_b S = -m^2 c^2$, introduced in Chapter 2. If you solve for S in the spacetime with the metric in Eq. (11.12), you can get the trajectories by the usual procedure of constructive interference. But we will follow a different procedure which will emphasize the power of the principle of equivalence.

*You can do it with
HJ ...*

To obtain this result, let us begin with the trajectory of a particle in special relativity under the action of a central force. The angular momentum $\mathbf{J} = \mathbf{r} \times \mathbf{p}$ is still conserved but the momentum is now given by $\mathbf{p} = \gamma m \mathbf{v}$ with $\gamma \equiv [1 - (v^2/c^2)]^{-1/2}$. So the relevant conserved component of the angular momentum is $J = mr^2(d\theta/d\tau) = \gamma mr^2(d\theta/dt)$ and not $mr^2(d\theta/dt)$. (This, incidentally, means that Kepler's second law regarding areal velocity does *not* hold in special relativistic motion in a central force in terms of the *coordinate time* t .) Consider now the motion of a *free* special relativistic particle described in *polar* coordinates. The standard

*... but the Principle
of equivalence is
more insightful*

relation $E^2 = \mathbf{p}^2 c^2 + m^2 c^4$ can be manipulated to give the equation

$$\frac{E^2}{c^2} \left(\frac{dr}{dt} \right)^2 = E^2 - \left(\frac{J^2 c^2}{r^2} + m^2 c^4 \right). \quad (11.15)$$

(This is still the description of a *free* particle moving in a straight line but in the polar coordinates!) Since special relativity must hold around any event, we can obtain the corresponding equation for general relativistic motion by simply replacing dr , dt by the proper quantities $\sqrt{|g_{11}|}dr$, $\sqrt{|g_{00}|}dt$ and the energy E by $E/\sqrt{|g_{00}|}$ (which is just the redshift obtained above) and $J = mr^2(d\theta/dt)$ to $J = mr^2(d\theta/\sqrt{|g_{00}|}dt)$ in this equation. This gives the equation for the orbit of a particle of mass m , energy E and angular momentum J around a body of mass M . With some simple manipulation, this can be written in a suggestive form as:

$$\left(1 - \frac{2GM}{c^2 r} \right)^{-1} \frac{dr}{dt} = \frac{c}{E} [E^2 - V_{\text{eff}}^2(r)]^{1/2} \quad (11.16)$$

with an effective potential:

$$V_{\text{eff}}^2(r) = m^2 c^4 \left(1 - \frac{2GM}{c^2 r} \right) \left(1 + \frac{J^2}{m^2 r^2 c^2} \right). \quad (11.17)$$

You can now work out various features of general relativistic orbits exactly as you do it in standard Kepler problem (see e.g., Chapter 25 of [50]). And the above derivation clearly shows that the particle is essentially following special relativistic, *free particle* motion at any event, in the locally inertial coordinates! (This is again a case of general relativity for the price of special relativity!)

For practical purposes, it is useful to rewrite Eq. (11.16) as a differential equation for $r(\theta)$. Noting that the expression for conserved angular momentum will also change from $J = mr^2(d\theta/dt)$ to $J = mr^2(d\theta/\sqrt{|g_{00}|}dt)$ and manipulating these equations, it is easy to obtain an expression for $dr/d\theta$. Differentiating this result will give the equation for the orbit in the standard form:

$$\frac{d^2 u}{d\theta^2} + u = \frac{GMm^2}{J^2} + \frac{3GM}{c^2} u^2. \quad (11.18)$$

The first term on the right hand side is purely Newtonian (see Eq. (3.42)) and the second term is the correction from general relativity. The ratio of these two terms is $(J/mrc)^2 \approx (v/c)^2$ where r and v are the typical radius and speed of the particle.

This correction term changes the nature of the orbits in two ways. First, it changes the relationship between the parameters of the orbit and the energy and angular momentum of the particle. More importantly, it makes

Try it out

Orbits in GR

the elliptical orbit of Newtonian gravity to precess slowly which is of greater observational importance. The exact solution to Eq. (11.18) can be given only in terms of elliptic functions and hence is not very useful. An approximate solution to Eq. (11.18), however, can be obtained fairly easily when the orbit has a very low eccentricity and is nearly circular (which is the case for most planetary orbits). Then the lowest order solution will be $u = u_0 = \text{constant}$ and one can find the next order correction by perturbations theory. This can be done *without* assuming that $2GMu_0/c^2 = 2GM/c^2 r_0$ is small, so that the result is valid *even for orbits close to the Schwarzschild radius*, as long as the orbit is nearly circular.

Let the radius of the circular orbit be r_0 for which $u = (1/r_0) \equiv k_0$. For the actual orbit, $u = k_0 + u_1$ where we expect the second term to be a small correction. Changing the variables from u to u_1 , where $u_1 = u - k_0$, Eq. (11.18) can be written as

$$u_1'' + u_1 + k_0 = \frac{GMm^2}{J^2} + \frac{3GM}{c^2} (u_1^2 + k_0^2 + 2u_1 k_0) . \quad (11.19)$$

We now choose k_0 to satisfy the condition

$$k_0 = \frac{GMm^2}{J^2} + k_0^2 \frac{3GM}{c^2} , \quad (11.20)$$

which determines the radius $r_0 = 1/k_0$ of the original circular orbit in terms of the other parameters. Now the equation for u_1 becomes

$$u_1'' + \left(1 - \frac{6k_0 GM}{c^2}\right) u_1 = \frac{3GM}{c^2} u_1^2 . \quad (11.21)$$

This equation is exact. We shall now use the fact that the deviation from circular orbit, characterized by u_1 is small and ignore the right hand side of equation Eq. (11.21). Solving Eq. (11.21), with the right hand side set to zero, we get

$$u_1 \cong A \cos \left[\left(1 - \frac{6GM}{c^2 r_0}\right)^{1/2} \theta \right] . \quad (11.22)$$

We see that r does not return to its original value at $\theta = 0$ when $\theta = 2\pi$ indicating a precession of the orbit. We encountered the same phenomenon in the case of motion in a Coulomb field as well in Chapter 3. As described in that context, the argument of the cosine function becomes 2π when

The precession, again!

$$\theta_c \approx 2\pi [1 - (6GM/c^2 r_0)]^{-1/2} , \quad (11.23)$$

which gives the precession $(\theta_c - 2\pi)$ per orbit. We can make a naive comparison between this precession rate and the corresponding one in the Coulomb problem by noticing that, in the latter case, we can substitute

$\alpha = GMm$ and $J^2 = GMm^2 r_0$ [which follows from Eq. (3.47)] to obtain

$$\omega_{\text{cl}}^2 \rightarrow 1 - \frac{GM}{c^2 r_0}, \quad (11.24)$$

*Quite different from
EM result*

which differs by a factor 6 in the corresponding term in general relativity.

If we attempt to reproduce the general relativistic results by an effective Newtonian potential, then a comparison with Eq. (3.42) tells us that we need to find a V_{eff} which satisfies the equation

$$-\frac{m}{J^2} \frac{dV_{\text{eff}}}{du} = \frac{GMm^2}{J^2} + \frac{3GM}{c^2} u^2, \quad (11.25)$$

which integrates to give

$$V_{\text{eff}} = -\frac{GMm}{r} - \frac{GMJ^2}{mc^2} \frac{1}{r^3}. \quad (11.26)$$

The trouble with this effective potential is that it depends on the angular momentum J of the particle which is somewhat difficult to motivate physically. But if you are willing to live with it, then one can introduce a pseudo Newtonian description of the general relativistic Kepler problem by taking the equations of motion to be $m(d^2 \mathbf{x}/d\tau^2) = \mathbf{F}$ with

A curiosity

$$\mathbf{F} = -\hat{\mathbf{r}} \frac{GMm}{r^2} \left(1 + \frac{3(\hat{\mathbf{r}} \times \mathbf{u})^2}{c^2} \right); \quad \mathbf{u} = \frac{d\mathbf{x}}{d\tau}, \quad (11.27)$$

where τ is the proper time. You can convince yourself that the conserved angular momentum now is $\mathbf{J} = m\mathbf{x} \times \mathbf{u}$ which will ensure that the above force reproduces the correct relativistic orbit equation. Unfortunately, this force law does not seem to lead to any other useful insight.

Classically, one thought of a black hole as a perfect absorber: Matter can fall into it but nothing can come out of it. In the early seventies, Bekenstein argued that this asymmetry can lead to the violation of second law of thermodynamics unless we associate an entropy with the black hole which is proportional to its area. This association made black holes rather peculiar thermodynamic objects. They were expected to possess an entropy and energy (given by Mc^2) but no temperature! This is because if black holes have a non-zero temperature, then they have to radiate a thermal spectrum of particles and this seemed to violate the classical notion that “nothing can come out of a black hole”. Given the fact that we do not know of any other system which possesses thermodynamic entropy and energy but not a temperature, this definitely looked peculiar.

System with energy and entropy but no temperature?

This puzzle was solved when, in the mid-seventies, Hawking discovered that a black hole *does* have a temperature, when viewed from a quantum mechanical perspective. A black hole which forms due to collapse of matter will emit — at late times — a thermal radiation which is characterized by this temperature. The rigorous derivation of this result requires a fair knowledge of quantum field theory but I will present, in this chapter, a simplified derivation which captures its essence [51, 52].

For a taste of history, see Box 12.1

We begin with a simple problem in special relativity but analyze it in a slightly unconventional way. Consider an inertial reference frame S and an observer who is moving at a speed v along the x -axis in this frame. If her trajectory is $x = vt$ then the clock she is carrying will show the proper time $\tau = t/\gamma$ where $\gamma = (1 - v^2/c^2)^{-1/2}$. Combining these results we can write her trajectory in parametrized form as

Doppler shift, from an unorthodox approach

$$t(\tau) = \gamma\tau; \quad x(\tau) = \gamma v\tau. \quad (12.1)$$

These equations give us her position in the spacetime when her clock reads τ .

Inertial motion

Suppose that a monochromatic plane wave, represented by the function $\phi(t, x) \equiv \exp -i\Omega(t - x/c)$, exists at all points in the inertial frame. This is clearly a plane wave of unit amplitude — as you will see soon, we don't care about the amplitude — and frequency Ω propagating along the positive x-axis. At any given x , it oscillates with time as $\exp(-i\Omega t)$ so Ω is the frequency as measured in S . The moving observer, of course, will measure how the ϕ changes with respect to *her* proper time. This is easily obtained by substituting the trajectory $t(\tau) = \gamma\tau; x(\tau) = \gamma v\tau$ into the function $\phi(t, x)$ obtaining $\phi[\tau] \equiv \phi[t(\tau), x(\tau)]$. A simple calculation gives

$$\begin{aligned}\phi[t(\tau), x(\tau)] &= \phi[\tau] = \exp[-i\tau\Omega\gamma(1 - v/c)] \\ &= \exp -i \left[\tau\Omega \sqrt{\frac{1 - v/c}{1 + v/c}} \right].\end{aligned}\quad (12.2)$$

Clearly, the observer sees a monochromatic wave with a frequency

$$\Omega' \equiv \Omega \sqrt{\frac{1 - v/c}{1 + v/c}}. \quad (12.3)$$

So an observer, moving with uniform velocity, will perceive a monochromatic wave as a monochromatic wave but with a Doppler shifted frequency; this is, of course, a standard result in special relativity derived in a slightly different manner.

Box 12.1: A little history

It all started with the theoretical discoveries in the seventies suggesting an intimate connection between thermodynamics and black holes with contributions from John Wheeler, Jacob Bekenstein, Stephan Hawking, Paul Davies, Bill Unruh and many others.

It occurred to John Wheeler that, by throwing a hot cup of tea into a black hole, he can hide the thermodynamic entropy of the tea forever from the observers who cannot access information from inside the black hole. This could allow a possible violation of second law of thermodynamics and Wheeler posed this problem to Bekenstein, who was at that time a graduate student [53]. Bekenstein came up with a remarkable solution to this difficulty.

Bekenstein suggested that the black hole should be associated with an entropy which is proportional to its area. When the cup of tea falls into the black hole, it increases the black hole's mass and size and hence the surface area. Bekenstein argued that — if the black holes have an entropy proportional to their area — then everything will be fine. In fact, Hawking had shown earlier that the areas of black holes

*Why black holes
must have entropy*

have a remarkable property: In any physical process involving normal matter and black holes or several black holes, the sum of the surface areas of the black holes can never decrease. This is very similar to the behaviour of entropy in thermodynamics giving credence to the idea of attributing an entropy to black holes which is proportional to the area.

But there was a serious problem with Bekenstein's idea which made several physicists, including Hawking, believe that this is just a mathematical analogy and that one cannot "really" attribute an entropy to the black hole. Since black holes have a mass, they certainly have an energy proportional to mass. If we now attribute an entropy proportional to the area (which will be proportional to the square of the mass), then one must also attribute a non-zero temperature to the black holes (which is inversely proportional to their mass). But if black holes have a non-zero temperature then they must radiate while the prevailing notion was that nothing comes out of black holes. Hence many physicists, originally refused to believe Bekenstein's idea and in fact Bekenstein had a hard time convincing others in the 1972 Les Houches summer school. (For a taste of history, see [54].)

Analogy or Truth?

But very soon, Hawking's own research showed that black holes do radiate, as though they have a non-zero temperature, thereby making everything consistent. Soon after, Paul Davies and Bill Unruh independently showed that the result is, in fact, far more general and occurs whenever a class of observers cannot receive information from certain region of spacetime. Black holes are described by just one kind of spacetimes in which this happens but the result is far more general. The fact that thermal radiation and temperature arises in these contexts illustrates yet again the power of mathematics which can tell us more than what we have originally assumed!

The real fun begins when we use the same procedure for a uniformly *accelerated* observer (sometimes called Rindler observer) along the x -axis. If we know the trajectory $t(\tau), x(\tau)$ of a uniformly accelerated observer, in terms of the proper time τ shown by the clock she carries, then we can determine $\phi[t(\tau), x(\tau)] = \phi[\tau]$ and repeat the previous analysis. So we first need to determine the trajectory $t(\tau), x(\tau)$ of a uniformly accelerated observer in terms of the proper time τ . Remembering that the equation of motion in special relativity is $d(m\gamma\mathbf{v})/dt = \mathbf{F}$, we can write the equation of motion for an observer moving with constant acceleration g along the x -axis as

Accelerated motion

$$\frac{d}{dt} \frac{v}{\sqrt{1-v^2/c^2}} = g. \quad (12.4)$$

This equation is trivial to integrate since g is a constant. Solving for $v = dx/dt$ and integrating once again, we can get the trajectory to be a hyperbola

$$x^2 - c^2 t^2 = c^4 / g^2, \quad (12.5)$$

with suitable choices for the initial conditions. We also know from special relativity that when a stationary clock registers a time interval dt , the moving clock will show a smaller proper time interval $d\tau = dt[1 - (v^2(t)/c^2)]^{1/2}$ where $v(t)$ is the instantaneous speed of the clock. Determining $v(t)$ from Eq. (12.5), one can find the relation between t and the proper time τ (as shown by a clock carried by the accelerated observer) as:

$$\tau = \int_0^t dt' \sqrt{1 - \frac{v^2(t')}{c^2}} = \frac{c}{g} \sinh^{-1} \left(\frac{gt}{c} \right). \quad (12.6)$$

Inverting this relation one can get t as a function of τ . Using Eq. (12.5) we can then express x in terms of τ and get the trajectory of the uniformly accelerated observer to be

$$x(\tau) = \frac{c^2}{g} \cosh \left(\frac{g\tau}{c} \right); \quad t(\tau) = \frac{c}{g} \sinh \left(\frac{g\tau}{c} \right). \quad (12.7)$$

This is exactly in the same spirit as the trajectory in Eq. (12.1) for an inertial observer except that we are now talking about a uniformly accelerated observer.

*Exponential
redshift*

We can now proceed exactly like in Eq. (12.2) to figure out how the *accelerated* observer views the monochromatic wave. We get:

$$\phi[t(\tau), x(\tau)] = \phi[\tau] = \exp i \frac{c}{g} [\Omega \exp - \frac{g\tau}{c}] = \exp i \theta(\tau). \quad (12.8)$$

Unlike in the case of uniform velocity, we now find that the phase $\theta(\tau)$ of the wave itself is decreasing exponentially with the proper time of the observer. Since the instantaneous frequency of the wave is the time derivative of the phase, $\omega(\tau) = -d\theta/d\tau$, we find that an accelerated observer will see the wave with an instantaneous frequency that is being exponentially redshifted:

$$\omega(\tau) = \Omega \exp \left(-\frac{g\tau}{c} \right). \quad (12.9)$$

*Calculate the power
spectrum ...*

Since this is not a monochromatic wave at all, the next best thing is to ask for the power spectrum of this wave which will tell us how it can be built out of monochromatic waves of different frequencies. (This is what an engineer would have done to analyse a time varying signal!) We will take the power spectrum of this wave to be $P(v) = |f(v)|^2$ where $f(v)$ is the Fourier transform of $\phi(t)$ with respect to t :

$$\phi(t) = \int_{-\infty}^{\infty} \frac{dv}{2\pi} f(v) e^{ivt}. \quad (12.10)$$

Evaluating this Fourier transform is an nice exercise in complex analysis and one can do it by changing to the variable $\Omega \exp[-(gt/c)] = z$ and analytically continuing to $\text{Im } z$. You will then find that:

$$f(\nu) = (c/g)(\Omega)^{-i\nu g/c} \Gamma(i\nu c/g) e^{-\pi\nu c/2g}, \quad (12.11)$$

where Γ is the standard Gamma function. Taking the modulus $|f(\nu)|^2$ using the identity $\Gamma(x)\Gamma(-x) = -\pi/x \sin(\pi x)$, we get

$$\nu |f(\nu)|^2 = \frac{\beta}{e^{\beta h\nu} - 1}; \quad \beta \equiv \frac{1}{k_B T} = \frac{2\pi c}{\hbar g}. \quad (12.12)$$

This leads to the remarkable result that the power, per logarithmic band in frequency, is a Planck spectrum with temperature $k_B T = (\hbar g/2\pi c)$! Also note that though $f(\nu)$ in Eq. (12.11) depends on Ω , the power spectrum $|f(\nu)|^2$ is independent of Ω . It does not matter what the frequency of the original wave was! The characteristic wavelength corresponding to this frequency is c^2/g ; its value is about 1 light year for earth's gravity — so the scope of experimental detection of this result is slim. (Incidentally, $c^2/g_{\text{earth}} \simeq 1$ light year gives a relation between earth's gravity and its orbital period around the sun; this is one of the cosmic coincidences which does not seem to have any deep significance.)

... to get something spectacular!

The moral of the story is simple: An exponentially redshifted complex wave will have a power spectrum which is thermal with a temperature proportional to the acceleration — which is responsible for the exponential redshift in the first place. This is the key to a quantum field theory result, due to Unruh, that a thermometer which is uniformly accelerated will behave as though it is immersed in a thermal bath.

There are two issues we have glossed over to get the correct result. First, we defined the Fourier transform in Eq. (12.10) with $e^{i\nu t}$, while the frequency of the original wave was $e^{-i\Omega t}$. So we are actually referring to the *negative* frequency component of a wave which has a positive frequency in the inertial frame. The second — and closely related issue — is that we have been working with complex wave modes, not just the real parts of them. Both these can be justified by a more rigorous analysis when these modes actually describe the vacuum fluctuations (see Chapter 19) in the inertial frame rather than some real wave. But the essential idea — and even the essential maths — is captured by this analysis.

Two remarks

So what about the temperature of black holes? Well, black holes produce an exponential redshift to the waves that propagate from close to the gravitational radius to infinity. To make the connection, we will recall two results from Chapter 11. First, the line element of a black hole is

Now for the black holes

$$ds^2 = - \left(1 - \frac{2GM}{c^2 r} \right) c^2 dt^2 + \left(1 - \frac{2GM}{c^2 r} \right)^{-1} dr^2 + r^2 (d\theta^2 + \sin^2 \theta d\phi^2). \quad (12.13)$$

Second, if $\omega(r)$ is the frequency of radiation emitted by a body of radius r and ω_∞ is the frequency with which this radiation is observed at large distances, then $\omega_\infty = \omega(r)(1 - 2GM/c^2r)^{1/2}$.

Let us now consider a wave packet of radiation emitted from a radial distance r_e at time t_e and observed at a large distance r at time t . The trajectory of the wave packet is, of course, given by $ds^2 = 0$ in Eq. (12.13) which — when we use $d\theta = d\phi = 0$ — is easy to integrate. (This result again follows from the principle of equivalence because, in the freely falling frame, light rays follow the trajectory with $ds^2 = 0$.) We get

$$\begin{aligned} c(t - t_e) &= r - r_e + \frac{2GM}{c^2} \ln \left(\frac{1 - 2GM/c^2r}{1 - 2GM/c^2r_e} \right) \\ &= r - r_e + \frac{4GM}{c^2} \ln \left(\frac{\omega_e}{\omega(r)} \right). \end{aligned} \quad (12.14)$$

For $r_e \gtrsim 2GM/c^2$, $r \gg 2GM/c^2$, this gives the frequency of radiation to be exponentially redshifted, as measured by an observer at infinity:

$$\omega(t) \propto \exp(-c^3t/4GM) \equiv K \exp(-(gt/c)), \quad (12.15)$$

where K is a constant (which turns out to be unimportant) and we have introduced the quantity

$$g = \frac{c^4}{4GM} = \frac{GM}{(2GM/c^2)^2}, \quad (12.16)$$

Again, the
exponential
redshift

which gives the gravitational acceleration GM/r^2 at the Schwarzschild radius $r = 2GM/c^2$ and is called the *surface gravity*. Once you have exponential redshift, the rest of the analysis proceeds as before. An observer detecting the exponentially redshifted radiation at late times ($t \rightarrow \infty$), originating from a region close to $r = 2GM/c^2$ will attribute to this radiation a Planckian power spectrum given by Eq. (12.12) which becomes:

$$k_B T = \frac{\hbar g}{2\pi c} = \frac{\hbar c^3}{8\pi GM}. \quad (12.17)$$

This forms the basis for associating a temperature with a black hole.

Once again, there is an extra (non-trivial) issue related to the question regarding the origin of the *complex* wave mode in the case of a black hole. The answer is the same as in the case of an accelerated observer we discussed earlier, with one interesting twist. Think of a spherical body surrounded by vacuum. In quantum theory, this vacuum will have a pattern of fluctuations which can be described in terms of complex wave modes. Suppose the body now collapses to form a black hole. The collapse upsets

the delicate balance between the wave modes in the vacuum and manifests — at late times — as thermal radiation propagating to infinity.

Using the expression in Eq. (12.17) for the temperature $T(M)$ of the black hole, and the energy of the black hole (Mc^2), we can formally integrate the relation $dS = dE/T$ to obtain the entropy of the black hole: *The entropy*

$$\frac{S}{k_B} = \int_0^M \frac{d(\bar{M}c^2)}{T(\bar{M})} = \pi \left(\frac{2GM}{c^2} \right)^2 \left(\frac{G\hbar}{c^3} \right)^{-1} = \frac{1}{4} \frac{4\pi r_H^2}{L_P^2}, \quad (12.18)$$

where $r_H = 2GM/c^2$ is the horizon radius of the black hole and $L_P = (G\hbar/c^3)^{1/2}$ is the so called Planck length. The entropy (which should be dimensionless when you use sensible units with $k_B = 1$) is just one quarter of the area of the horizon in units of Planck length. Getting this factor $1/4$ is a holy grail in models for quantum gravity — but that is another story.

Box 12.2: The thermodynamics behind Einstein's equations

There is a remarkable connection between the first law of thermodynamics and the laws describing gravity in a wide class of theories. The figure below illustrates this analogy. The figure on the left shows some amount of gas confined to a box, the volume of which can be changed by moving the piston. If you let the piston move outward due to the pressure of the gas, one can extract some mechanical work from the gas as well as change the internal energy of the gas. The gas can also exchange heat with the surroundings that can be expressed in terms of the entropy change of the gas. The first law of thermodynamics relates the changes in these three quantities: entropy, internal energy and mechanical work.

Let us move on from a box of gas to a spacetime with a horizon (see the figure on the right). The location of the horizon in this figure plays a role analogous to the position of the piston in the figure on the left side. While you cannot push around a horizon, you can certainly consider two different spacetimes with the horizons at two slightly different locations. This displacement of the horizon again causes changes in the properties of the spacetime which — in turn — are governed by the equation describing the gravity. Remarkably enough, one can prove that this equation reduces to a form identical to the equation in the case of gas with a piston!

Deep, not completely understood

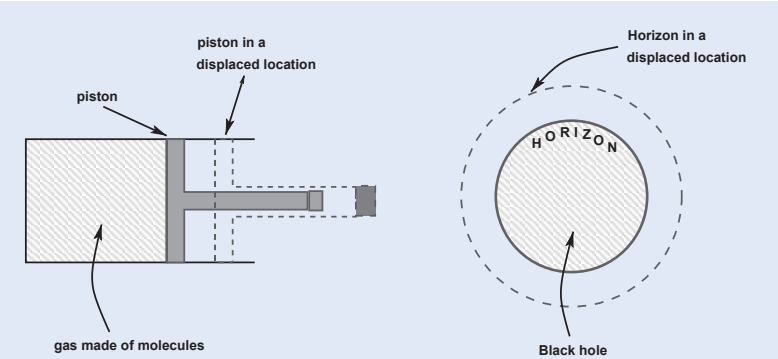


Fig. 12.1: Analogy between Einstein’s equations and thermodynamics. See text for discussion.

This result was first obtained by me in 2002 in the simplest context of horizons in Einstein’s theory [55, 56]. Further work by different groups has now established that this result is true for a wide class of theories of gravity much more general than Einstein’s theory. In the general context, the temperature associated with the horizon is independent of the theory one is studying but the entropy depends crucially on the theory. Remarkably enough, the thermodynamic description picks out the correct expression for entropy in each theory thereby showing that the entropy density associated with a horizon contains the necessary information to reconstruct the underlying theory. This — and several other results — suggest that gravity is an emergent phenomenon like e.g elasticity or fluid mechanics and the field equations of gravity only have the same status as the equations of fluid mechanics. This emergent gravity paradigm — which is a major research area today — is a direct offshoot of the results discussed in this chapter!

Consider an electron (with a spin) orbiting in an atom, treated along classical lines. In the instantaneous rest frame of the orbiting electron, the Coulomb field (Ze^2/r^2) of the nucleus gives rise to a magnetic field $(v/c)(Ze^2/r^2)$. This magnetic field couples to the magnetic moment ($e\hbar/2m_e c$) of the electron, thereby contributing to the effective energy of coupling between the spin and orbital motion. Clearly, this is a special relativistic effect of the order of $(v/c)^2$. But if you compare this naive theoretical result with observation, you will find that they differ by a factor $(1/2)$. This factor is also due to a relativistic effect called the Thomas precession. It is one of the peculiar features of special relativity which is purely kinematic in origin and has observational consequences [57].

One practical reason why this result is important

This precession also has an interesting geometrical interpretation which allows one to relate it to another — apparently unconnected — physical phenomenon, viz. the rotation of the plane of the Foucault pendulum. In this chapter, I will provide a straightforward (and possibly not very inspiring) derivation of the Thomas precession. In the next chapter we will explore the Foucault pendulum and the geometrical relationship between the two.

Consider the standard Lorentz transformation equations between two inertial frames which are in relative motion along the x -axis with a speed $V \equiv c\beta$. These are given by $x = \gamma(x' + Vt')$, $t = \gamma(t' + Vx'/c^2)$ where $\gamma = (1 - \beta^2)^{-1/2}$. We know that the quantity $s^2 \equiv (-c^2t^2 + |\mathbf{x}|^2)$ remains invariant under the Lorentz transformation. A quadratic expression of this form is similar to the length of a vector in three dimensions which is invariant under rotation of the coordinate axes. This suggests that the transformation between the inertial frames can be thought of as a rotation in four dimensional space. The rotation must be in the $t - x$ plane characterized by a parameter, say, ψ . Indeed, the Lorentz transformation can be equivalently written as

Lorentz transformation = rotation by imaginary angle

$$x = x' \cosh \psi + ct' \sinh \psi, \quad ct = x' \sinh \psi + ct' \cosh \psi, \quad (13.1)$$

with $\tanh \psi = (V/c)$, which determines the parameter ψ (called the *rapidity*) in terms of the relative velocity between the two frames. Eq. (13.1) can be thought of as a rotation by a complex angle $i\psi$.

Two successive Lorentz transformations with velocities V_1 and V_2 , *along the same direction* x , correspond to two successive rotations in the t - x plane by angles, say, ψ_1 and ψ_2 . Since two rotations in the same plane commute, it is obvious that these two Lorentz transformations commute and are equivalent to a rotation by an angle $\psi_1 + \psi_2$ in the t - x plane. This results in a single Lorentz transformation with a velocity parameter given by the relativistic sum of the two velocities V_1 and V_2 . Note that the rapidities simply add while the velocity addition formula is more complicated.

The situation, however, changes in the case of Lorentz transformations along two different directions. This will correspond to rotations in two different planes and it is well known that such rotations will not commute. The order in which the Lorentz transformations are carried out is important if they are along different directions. Suppose a frame S_1 is moving with a velocity $\mathbf{V}_1 = V_1 \mathbf{n}_1$ (where \mathbf{n}_1 is a unit vector) with respect to a reference frame S_0 and we do a Lorentz boost to connect the coordinates of these two frames. Now suppose we do another Lorentz boost with a velocity $\mathbf{V}_2 = V_2 \mathbf{n}_2$ to go from S_1 to S_2 . We want to know what kind of transformation will now take us directly from S_0 to S_2 . If $\mathbf{n}_1 = \mathbf{n}_2$, then the two Lorentz transformations are along the same axis and one can go from S_0 to S_2 by a single Lorentz transformation. But this is not possible if the two directions \mathbf{n}_1 and \mathbf{n}_2 are different. It turns out that, in addition to the Lorentz transformation, one also has to rotate the spatial coordinates by a particular amount.

This is the root cause of Thomas precession. For a body moving in an accelerated trajectory with the direction of velocity vector changing continuously, the instantaneous Lorentz frames are obtained by boosts along different directions at each instant. Since such successive boosts are equivalent to a boost plus a rotation of spatial axes, there is an effective rotation of the coordinate axes which occurs in the process. If the body carries an intrinsic vector (like spin) with it, the orientation of that vector will undergo a shift.

After all that English, let us to establish the idea mathematically. To do this, we need the Lorentz transformations connecting two different frames of reference, when one of them is moving along an arbitrary direction \mathbf{n} with speed $V \equiv \beta c$. The time coordinates are related by the obvious formula

$$x^0 = \gamma(x^0 - \boldsymbol{\beta} \cdot \mathbf{x}) , \quad (13.2)$$

where we are using the notation $x^i = (x^0, \mathbf{x}) = (ct, \mathbf{x})$ to denote the four-vector coordinates. To obtain the transformation of the spatial coordinate, we first write the spatial vector \mathbf{x} as a sum of two vectors; $\mathbf{x}_{\parallel} = \mathbf{V}(\mathbf{V} \cdot \mathbf{x})/V^2$

Two Lorentz transformations = One Lorentz transformation + Rotation

The origin of Thomas precession

Combining two Lorentz transformations

which is parallel to the velocity vector and $\mathbf{x}_\perp = \mathbf{x} - \mathbf{x}_\parallel$ which is perpendicular to the velocity vector. We know that, under the Lorentz transformation, we have $\mathbf{x}'_\perp = \mathbf{x}_\perp$ while $\mathbf{x}'_\parallel = \gamma(\mathbf{x}_\parallel - \mathbf{V}t)$. Expressing everything in terms of \mathbf{x} and \mathbf{x}' , it is easy to show that the final result can be written in the vectorial form (with $\boldsymbol{\beta} = \beta \mathbf{n}$) as:

$$\mathbf{x}' = \mathbf{x} + \frac{(\gamma - 1)}{\beta^2} (\boldsymbol{\beta} \cdot \mathbf{x}) \boldsymbol{\beta} - \gamma \boldsymbol{\beta} x^0. \quad (13.3)$$

Equations (13.2) and (13.3) give the Lorentz transformation between two frames whose relative direction of motion is arbitrary.

We will now use this result to determine the effect of two consecutive Lorentz transformations for the case in which both $\mathbf{V}_1 = V_1 \mathbf{n}_1$ and $\mathbf{V}_2 = V_2 \mathbf{n}_2$ are small in the sense that $V_1 \ll c$, $V_2 \ll c$. Let the first Lorentz transformation take the four vector $x^b = (ct, \mathbf{x})$ to x^b_1 and the second Lorentz transformation take this further to x^a_{21} . Performing the same two Lorentz transformations in reverse order leads to the vector which we will denote by x^a_{12} . We are interested in the difference $\delta x^a \equiv x^a_{21} - x^a_{12}$ to the lowest non-trivial order in (V/c) . Since this involves the product of two Lorentz transformations, we need to compute it keeping all terms up to *quadratic* order in V_1 and V_2 . Explicit computation, using, Eq. (13.2) and Eq. (13.3) now gives

Try it out!

$$\begin{aligned} x^0_{21} &\approx [1 + \frac{1}{2}(\boldsymbol{\beta}_2 + \boldsymbol{\beta}_1)^2]x^0 - (\boldsymbol{\beta}_2 + \boldsymbol{\beta}_1) \cdot \mathbf{x} \\ \mathbf{x}_{21} &\approx \mathbf{x} - (\boldsymbol{\beta}_2 + \boldsymbol{\beta}_1)x^0 + [\boldsymbol{\beta}_2(\boldsymbol{\beta}_2 \cdot \mathbf{x}) + \boldsymbol{\beta}_1(\boldsymbol{\beta}_1 \cdot \mathbf{x})] + \boldsymbol{\beta}_2(\boldsymbol{\beta}_1 \cdot \mathbf{x}), \end{aligned} \quad (13.4)$$

accurate to $\mathcal{O}(\beta^2)$. It is obvious that terms which are symmetric under the exchange of 1 and 2 in the above expression will cancel out when we compute $\delta x^a \equiv x^a_{21} - x^a_{12}$. Hence, we get $\delta x^0 = 0$ to this order of accuracy. In the spatial components, the only surviving term is the one arising from last term in the expression for \mathbf{x}_{21} , which gives

$$\delta \mathbf{x} = [\boldsymbol{\beta}_2(\boldsymbol{\beta}_1 \cdot \mathbf{x}) - \boldsymbol{\beta}_1(\boldsymbol{\beta}_2 \cdot \mathbf{x})] = \frac{1}{c^2} (\mathbf{V}_1 \times \mathbf{V}_2) \times \mathbf{x}. \quad (13.5)$$

Comparing this with the standard result for infinitesimal rotation of coordinates, $\delta \mathbf{x} = \boldsymbol{\Omega} \times \mathbf{x}$, we find that the net effect of two Lorentz transformations leaves a residual *spatial rotation* about the direction $\mathbf{V}_1 \times \mathbf{V}_2$. Since this result is obtained by taking the difference between two successive Lorentz transformations, $\delta \mathbf{x} \equiv \mathbf{x}_{21} - \mathbf{x}_{12}$, we can think of each one contributing an effective rotation by the amount $(1/2)(\mathbf{V}_1 \times \mathbf{V}_2)/c^2$.

Consider now a particle with a spin moving in a circular orbit. (For example, it could be an electron in an atom; the classical analysis continues to apply essentially because the effect is purely kinematic!). At two instances in time t and $t + \delta t$, the velocity of the electron will be in different

We can get away with classical analysis

directions \mathbf{V}_1 and $\mathbf{V}_1 + \mathbf{a}\delta t$ where \mathbf{a} is the acceleration. This should lead to a change in the angle of orientation of the axes by the amount

$$\delta\boldsymbol{\Omega} = \frac{1}{2} \frac{(\mathbf{V}_1 \times \mathbf{V}_2)}{c^2} = \frac{1}{2} \frac{(\mathbf{V}_1 \times \mathbf{a})}{c^2} \delta t, \quad (13.6)$$

corresponding to the angular velocity

$$\boldsymbol{\omega} = \frac{\delta\boldsymbol{\Omega}}{\delta t} = \frac{1}{2} \frac{\mathbf{V}_1 \times \mathbf{a}}{c^2}. \quad (13.7)$$

This is indeed the correct expression for Thomas precession in the non-relativistic limit (since we had assumed $V_1 \ll c, V_2 \ll c$).

Let me now outline a rigorous derivation of this effect which is valid for even relativistic speeds. To set the stage, we again begin with the rotations in 3-dimensional space. A given rotation can be defined by specifying the unit vector \mathbf{n} in the direction of the axis of rotation and the angle θ through which the axes are rotated. We can associate with this rotation a 2×2 matrix

$$R(\theta) = \cos(\theta/2) - i(\boldsymbol{\sigma} \cdot \mathbf{n}) \sin(\theta/2) = \exp - \left[\frac{i\theta}{2} (\boldsymbol{\sigma} \cdot \mathbf{n}) \right], \quad (13.8)$$

where $\boldsymbol{\sigma}_\alpha$ are the standard Pauli matrices and the $\cos(\theta/2)$ term is considered to be multiplied by the unit matrix though it is not explicitly indicated. The equivalence of the two forms — the exponential and trigonometric — of $R(\theta)$ in Eq. (13.8) can be demonstrated by expanding the exponential in a power series and using the easily proved relation $(\boldsymbol{\sigma} \cdot \mathbf{n})^2 = 1$. We can also associate with a 3-vector \mathbf{x} the 2×2 matrix $X = \mathbf{x} \cdot \boldsymbol{\sigma}$. The effect of any rotation can now be concisely described by the matrix relation $X' = RXR^*$.

Since we can think of Lorentz transformations as rotations by an imaginary angle, all these results generalize in a natural way to the Lorentz transformations. We can associate with a Lorentz transformation in the direction \mathbf{n} with the speed $V = c \tanh \alpha$, the 2×2 matrix

$$L = \cosh(\alpha/2) + (\mathbf{n} \cdot \boldsymbol{\sigma}) \sinh(\alpha/2) = \exp \frac{1}{2} (\boldsymbol{\alpha} \cdot \boldsymbol{\sigma}). \quad (13.9)$$

The change from trigonometric functions to hyperbolic functions is in accordance with the fact that Lorentz transformations correspond to rotation by an *imaginary* angle. Just as in the case of rotations, we can associate to any event $x^i = (x^0, \mathbf{x})$ a (2×2) matrix $P \equiv x^i \sigma_i$ where σ_0 is the identity matrix and σ_α are the Pauli matrices. Under a Lorentz transformation along the direction \mathbf{n} with speed V , the event x^i goes to x'^i and P goes P' . (By convention, the σ_i 's do not change.) They are related by

$$P' = LPL^*, \quad (13.10)$$

where L is given by Eq. (13.9).

Warm up: 3D rotations

Why $\theta/2$?

A rotation through an angle θ about a given axis is due to successive reflections in two planes which meet along the axis at an angle $\theta/2$.

Now for the real thing; let $3 \rightarrow 4$

Consider an inertial, laboratory frame S_0 and let $S(t)$ be a Lorentz frame co-moving with a particle (which has a non-zero spin) at time t . These two frames are related to each other by a Lorentz transformation with a velocity \mathbf{V} . Consider a pure Lorentz boost in the *comoving* frame of the particle which changes its velocity relative to the lab frame from \mathbf{V} to $\mathbf{V} + d\mathbf{V}$. We know that the resulting final configuration cannot be reached from S_0 by a pure boost and we require a rotation by some angle $\delta\boldsymbol{\theta} = \boldsymbol{\omega}dt$ followed by a simple boost. This leads to the relation, in terms of the 2×2 matrices corresponding to the rotation and Lorentz transformations, as:

Kinematics: defined precisely

$$L(\mathbf{V} + d\mathbf{V})R(\boldsymbol{\omega}dt) = L_{\text{comov}}(d\mathbf{V})L(\mathbf{V}) . \quad (13.11)$$

The right hand side represents, in matrix form, two Lorentz transformations. The left hand side represents the same effect in terms of one Lorentz transformation and one rotation — the parameters of which are at present unknown. In the right hand side of Eq. (13.11), the matrix $L_{\text{comov}}(d\mathbf{V})$ has a subscript “comoving” to stress the fact that this operation corresponds to a pure boost *only* in the comoving frame and *not* in the lab frame. To take care of this, we do the following: We first bring the particle to rest by applying the inverse Lorentz transformation operator $L^{-1}(\mathbf{V}) = L(-\mathbf{V})$. Then we apply a boost $L(\mathbf{a}_{\text{comov}}d\tau)$ where $\mathbf{a}_{\text{comov}}$ is the acceleration of the system in the comoving frame. Since the object was at rest initially, *this second operation* can be characterized by a pure boost. Finally, we transform back from the lab to the moving frame by applying $L(\mathbf{V})$. We thus obtain the relation

Note the meaning of ‘comoving’

$$L_{\text{comov}}(d\mathbf{V}) = L(\mathbf{V})L(\mathbf{a}_{\text{comov}}d\tau)L(-\mathbf{V}) . \quad (13.12)$$

Using this in Eq. (13.11), we get $L(\mathbf{V} + d\mathbf{V})R(\boldsymbol{\omega}dt) = L(\mathbf{V})L(\mathbf{a}_{\text{comov}}d\tau)$. In this equation, the unknowns are $\boldsymbol{\omega}$ and $\mathbf{a}_{\text{comov}}$. Moving the unknown terms to the left hand side, we have the equation,

$$R(\boldsymbol{\omega}dt)L(-\mathbf{a}_{\text{comov}}d\tau) = L(-[\mathbf{V} + d\mathbf{V}])L(\mathbf{V}) , \quad (13.13)$$

which can be solved for $\boldsymbol{\omega}$ and $\mathbf{a}_{\text{comov}}$. If we denote the rapidity parameters for the two infinitesimally separated Lorentz boosts by α and $\alpha' \equiv \alpha + d\alpha$, and the corresponding directions by \mathbf{n} and $\mathbf{n}' \equiv \mathbf{n} + d\mathbf{n}$, then this matrix equation can be expanded to first order quantities to give

$$1 - (i\boldsymbol{\omega}dt + \mathbf{a}d\tau) \cdot \frac{\boldsymbol{\sigma}}{2} = [\cosh(\alpha'/2) - (\mathbf{n}' \cdot \boldsymbol{\sigma}) \sinh(\alpha'/2)] \\ \times [\cosh(\alpha/2) - (\mathbf{n} \cdot \boldsymbol{\sigma}) \sinh(\alpha/2)] . \quad (13.14)$$

Performing the necessary Taylor series expansion in $d\alpha$ and $d\mathbf{n}$ in the right hand side and identifying the corresponding terms on both sides, we

find that $\mathbf{a}_{\text{comov}} = \mathbf{n}(d\alpha/d\tau) + (\sinh \alpha)(d\mathbf{n}/d\tau)$, and more importantly,

$$\boldsymbol{\omega} = (\cosh \alpha - 1) \left(\frac{d\mathbf{n}}{dt} \times \mathbf{n} \right), \quad (13.15)$$

The result for $\boldsymbol{\omega}dt$ has a nice geometrical interpretation; see next chapter

with $\tanh \alpha = (V/c)$. Expressing everything in terms of the velocity, it is easy to show that the expression for $\boldsymbol{\omega}$ is equivalent to

$$\boldsymbol{\omega} = \frac{\gamma^2}{\gamma+1} \frac{\mathbf{a} \times \mathbf{V}}{c^2} = (\gamma-1) \frac{(\mathbf{V} \times \mathbf{a})}{V^2}. \quad (13.16)$$

In the non-relativistic limit, this gives a precessional angular velocity $\boldsymbol{\omega} \cong (1/2c^2)(\mathbf{a} \times \mathbf{V})$ which the spin will undergo because of the non-commutativity of Lorentz transformations in different directions.

After all that, a simple derivation!

Having provided a fairly rigorous derivation of this effect, we will now describe a simple intuitive way of understanding the same [58]. This involves interpreting the extra rotation which arises when successive Lorentz transformations are performed in terms of the length contraction. Consider an aircraft flying around in a large circular orbit which we approximate by a polygon of N sides — with the understanding that, eventually, we will take the $N \rightarrow \infty$ limit. Once the aircraft traverses the N -gon, it is back to the starting point. In the laboratory frame, it has rotated through an angle 2π , but — in the airplane's instantaneous frame — traversing each side of the N -gon leads to a different result (see Fig. 13.1). While turning through an angle θ , the transverse distance is

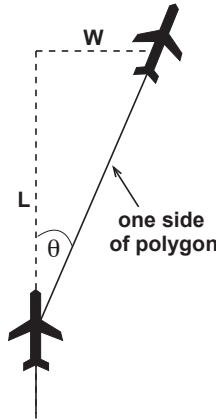


Fig. 13.1: An intuitive interpretation of the Thomas precession. We approximate a circular orbit as one made of a polygon with very large number of sides. While turning from one side to another side of the polygon, the transverse and longitudinal length scales transform differently with respect to the co-moving Lorentz frames. This effect accumulates to give the standard result for the Thomas precession when the orbit is completed.

still W but the longitudinal distance undergoes Lorentz contraction to become L/γ . Therefore, the angle of turn experienced at each vertex may be thought of as W divided by L/γ , giving $\gamma\theta$. So, the net effect is that, over a round trip the airplane has rotated with respect to local inertial frames by an amount $2\pi\gamma$ while it has rotated through 2π with respect to the laboratory frame. So the net *extra* rotation over the circular trip, completed in time T , say, is $\Delta\theta = 2\pi(\gamma - 1)$. The effective precession rate will then be

$$\frac{\omega_P}{\omega} \equiv \frac{\Delta\theta/T}{2\pi/T} = \gamma - 1. \quad (13.17)$$

This is same as Eq. (13.6) we obtained earlier, for the case of circular motion. While the argument is not rigorous, it certainly provides an intuitive understanding of what a bunch of Lorentz transformations can do.

Box 13.1: Geometrical way of combining rotations

An interesting issue in the study of rotations in 3-dimensional space is to characterize geometrically the effect of combining two arbitrary rotations [50]. You might enjoy proving the following construction for finding the resultant of two spatial rotations characterized by the directions $\mathbf{n}_1, \mathbf{n}_2$ and angles θ_1, θ_2 .

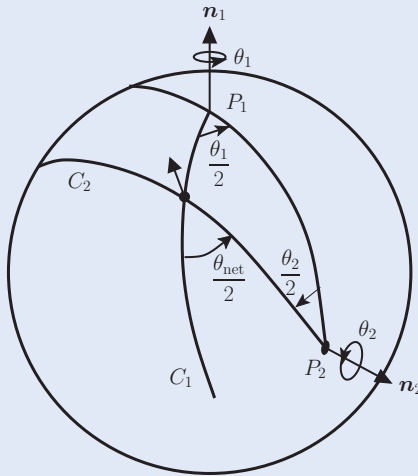


Fig. 13.2: A geometrical way of combining two rotations around two arbitrary axes.

Let the directions \mathbf{n}_1 and \mathbf{n}_2 be denoted by the points P_1 and P_2 on the surface of a unit sphere. Draw the great circle going through P_1 and P_2 . Draw another great circle C_1 passing through P_1 making an

angle $\theta_1/2$ with the circle P_1P_2 , i.e., the tangents to the two circles drawn at P_1 make an angle $\theta_1/2$. Similarly, draw a great circle C_2 passing through P_2 making an angle $\theta_2/2$ with the circle P_1P_2 . The orientations are to be as indicated in [Fig. 13.2](#). The intersection of C_1 and C_2 will give the direction of the axis of the resultant rotation and the external angle of the spherical triangle at the intersection will give $\theta/2$ where θ is the resultant angle of rotation.

The Pantheon in Paris was used by Leon Foucault on 31 March 1851, under the reign of Louis-Napoleon Bonaparte, the first titular president of the French republic, to give an impressive demonstration. Using a pendulum (with a 67 meter wire and a 28 kg pendulum bob), he could demonstrate the rotation of the Earth in a tell-tale manner. As the pendulum kept swinging, one could see that the plane of oscillation of the pendulum itself was rotating in a clockwise direction (when viewed from the top). The frequency of this rotation was $\omega = \Omega \cos \theta$ where Ω is the angular velocity of Earth and θ is the co-latitude of Paris. (That is, θ is the standard polar angle in spherical polar coordinates with the z -axis being the axis of rotation of Earth. So $\pi/2 - \theta$ is the geographical latitude). Foucault claimed, quite correctly, that this demonstrates the rotation of the Earth using an ‘in situ’ experiment without us having to look at the celestial objects.

Foucault’s demonstration

This result is quite easy to understand if the experiment was performed at the poles or the equator (instead of at Paris!). The situation at the north pole is as shown in Fig. 14.1. Here we see the Earth as rotating (from west to east, in the counter-clockwise direction when viewed from the top) underneath the pendulum, making one full turn in 24 hours. It appears reasonable to deduce from this that, as viewed from Earth, the plane of oscillation of the pendulum will make one full rotation in 24 hours; so the angular frequency ω of the rotation of the plane of the Foucault pendulum is just $\omega = \Omega$. (Throughout the discussion it is the rotation of the *plane of oscillation* of the pendulum we are concerned with; not the period of the pendulum $2\pi/\nu$, which — of course — is given by the standard formula involving the length of the suspension wire, etc.). At the equator, on the other hand, the plane of oscillation does not rotate. So the formula, $\omega = \Omega \cos \theta$, captures both limits correctly.

If only Paris was at the North Pole ...

It is easy to write down the equations of motion for the pendulum bob in the rotating frame of the Earth and solve them to obtain this result [36, 59] correct to linear order in Ω . Essentially, the Foucault pendulum effect arises due to the Coriolis force in the rotating frame of the Earth which

The most efficient (but unimaginative) way of deriving the result

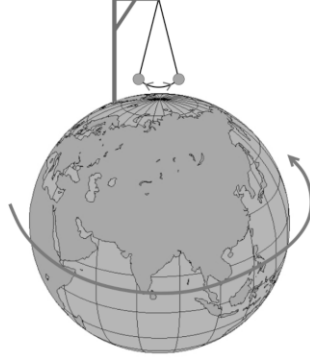


Fig. 14.1: The rotation of the plane of the Foucault pendulum is easy to understand if the pendulum was located in the north pole. However, as discussed in the text, the apparent simplicity of this result is deceptive.

leads to an acceleration $2\mathbf{v} \times \boldsymbol{\Omega}$ where \mathbf{v} , the velocity of the pendulum bob, is directed tangential to the Earth's surface to a good approximation. If we choose a local coordinate system with the Z -axis pointing normal to the surface of the Earth and the X, Y coordinates in the tangent plane at the location, then it is easy to show that the equations of motion for the pendulum bob are well approximated by

$$\ddot{X} + v^2 X = 2\Omega_z \dot{Y}; \quad \ddot{Y} + v^2 Y = -2\Omega_z \dot{X}, \quad (14.1)$$

where v is the period of oscillation of the pendulum and $\Omega_z = \Omega \cos \theta$ is the normal component of Earth's angular velocity. (This can be easily derived from Eq. (6.3) of Chapter 6.) In arriving at these equations we have ignored terms quadratic in Ω^2 and the vertical displacement of the pendulum. The solution to this equation is obtained by introducing the variable $q(t) \equiv X(t) + iY(t)$, which satisfies the equation

$$\ddot{q} + 2i\Omega_z \dot{q} + v^2 q = 0. \quad (14.2)$$

The solution, to the order of accuracy we are working with, is given by

$$q = X(t) + iY(t) = (X_0(t) + iY_0(t)) \exp(-i\Omega_z t), \quad (14.3)$$

where $X_0(t), Y_0(t)$ is the trajectory of the pendulum in the absence of Earth's rotation. It is clear that the net effect of rotation is to cause a shift in the plane of rotation at the rate $\Omega_z = \Omega \cos \theta$. Based on this knowledge and the results for the pole and the equator one can give a 'pure English' derivation of the result for intermediate latitudes by saying something like: "Obviously, it is the component of $\boldsymbol{\Omega}$ normal to the Earth at the location of the pendulum which matters and hence $\omega = \Omega \cos \theta$."

The first-principle approach, based on Eq. (14.1), of course, has the advantage of being rigorous and algorithmic; for example, if the effects of the ellipticity of the Earth are to be incorporated, it can be done by working with the equations of motion. But it does not give you an intuitive understanding of what is going on, and much less a unified view of this problem with other related problems having the same structure. We shall now describe an approach to this problem which has the advantage of providing a clear geometrical picture and connecting it up — somewhat quite surprisingly — with Thomas precession discussed in the previous chapter [60].

But, what is actually happening?

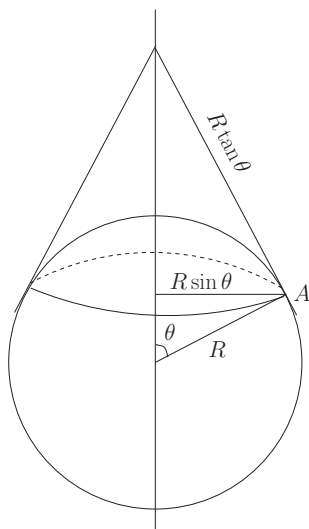


Fig. 14.2: A Foucault pendulum is located at the co-latitude θ (i.e., at the geographical latitude $(\pi/2 - \theta)$). A cone which is tangential at this latitude allows us to obtain a geometrical interpretation of the rotation of the plane of the Foucault pendulum.

An issue that causes some confusion as regards the Foucault pendulum is the following. While analyzing the behavior of the pendulum at the pole, one assumes that the plane of rotation remains fixed while the Earth rotates underneath it. If we make the same claim for a pendulum experiment done at an intermediate latitude, — i.e., if we say that the plane of oscillation remains invariant with respect to, say, the “fixed stars” and the Earth rotates underneath it — it seems natural that the period of rotation of the pendulum plane should *always* be 24 hours irrespective of the location! This, of course, is not true and it is also intuitively obvious that nothing happens to the plane of rotation at the equator. In this way of approaching the problem, it is not very clear how exactly the Earth’s rotation influences the motion of the pendulum.

A minor paradox, usually glossed over

The geometrical insight: parallel transport

We will now provide a geometrical approach to this problem, by rephrasing it as follows [61, 62]. The plane of oscillation of the pendulum can be characterized either by a vector normal to its plane or — equivalently — by a vector which is lying in the plane *and* tangential to the Earth’s surface. Let us now introduce a cone which is coaxial with the axis of rotation of the Earth and having its surface tangential to the Earth at the latitude of the pendulum (see Fig. 14.2). The base radius of such a cone will be $R \sin \theta$ where R is the radius of the Earth and the slant height of the cone will be $R \tan \theta$. Such a cone can be built out of a sector of a circle (as shown in Fig. 14.3) having the circumference $2\pi R \sin \theta$ and radius $R \tan \theta$ by identifying the lines OA and OB . The ‘deficit angles’ of the cone, α and $\beta \equiv 2\pi - \alpha$, satisfy the relations:

$$(2\pi - \alpha)R \tan \theta = 2\pi R \sin \theta , \tag{14.4}$$

which gives

$$\alpha = 2\pi(1 - \cos \theta); \qquad \beta = 2\pi \cos \theta . \tag{14.5}$$

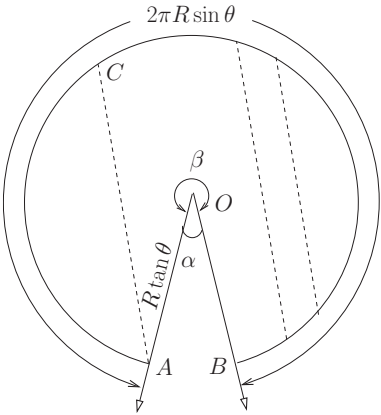


Fig. 14.3: The cone in the previous figure is built from the sector of the circle shown here. The parallel transport of a vector is easier to understand in terms of the deficit angle of the sector.

Analyse it on the cone

The behavior of the plane of the Foucault pendulum can be understood very easily in terms of this cone. Initially, the Foucault pendulum starts oscillating in some arbitrary direction at the point A, say. This direction of oscillation can be indicated by some straight line drawn along the surface of the cone (like AC in Fig. 14.3). While the plane of oscillation of the pendulum will rotate with respect to a coordinate system fixed on the Earth, it will always coincide with the lines drawn on the cone which remain fixed relative to the fixed stars. When the Earth makes one rotation,

we move from A to B in the flattened out cone in Fig. 14.3. Physically, of course, we identify the two points A and B with the same location on the surface of the Earth. But when a vector has been moved around a curve along the lines described above, on the curved surface of Earth, its orientation does not return to the original value. It is obvious from Fig. 14.3 that the orientation of the plane of rotation (indicated by a vector in the plane of rotation and tangential to the Earth's surface at B) is different from the corresponding vector at A . This process is called parallel transport and the fact that a vector changes on parallel transport around an arbitrary closed curve on a curved surface is a well known result in differential geometry and general relativity.

Clearly, the orientation of the vector changes by an angle $\beta = 2\pi \cos \theta$ during one rotation of Earth with period T . Since the rate of change is uniform throughout because of the steady state nature of the problem, the angular velocity of the rotation of the pendulum plane is given by

$$\omega = \frac{\beta}{T} = \frac{2\pi}{T} \cos \theta = \Omega \cos \theta . \quad (14.6)$$

This is precisely the result we were after. The key geometrical idea was to relate the rotation of the plane of the Foucault pendulum to the parallel transport of a vector characterizing the plane, around a closed curve on the surface of Earth. When this closed curve is not a geodesic — and we know that a curve of constant latitude is not a geodesic — the orientation of this vector changes when it completes one loop. There are sophisticated ways of calculating how much the orientation changes for a given curve on a curved surface. But in the case of a sphere, the trick of an enveloping cone provides a simple procedure. When the pendulum is located at the equator, the closed curve is the equator itself. The equator, being a great circle, is a geodesic on the sphere and hence the vector does not get ‘disoriented’ on going around it. So the plane of the pendulum does not rotate in this case. (In fact, there is a nice relation between the area enclosed by a curve on the sphere and the amount of rotation the vector will undergo when parallel transported around this curve; see the Appendix to this chapter.)

There you are!

Remarkably enough, one can show that an almost identical approach allows one to determine the Thomas precession of the spin of a particle (say, an electron) moving in a circular orbit around a nucleus [63].

This is good, but it gets better!

We saw in the last chapter that the rate of Thomas precession is given, in general, by an expression of the form

$$\boldsymbol{\omega} dt = (\cosh \chi - 1) (d\hat{\mathbf{n}} \times \hat{\mathbf{n}}) , \quad (14.7)$$

where $\tanh \chi = v/c$ and v is the velocity of the particle. In the case of a particle moving on a circular trajectory, the magnitude of the velocity remains constant and we can integrate this expression to obtain the net angle of precession during one orbit. For a circular orbit, $d\hat{\mathbf{n}}$ is always

perpendicular to $\hat{\mathbf{n}}$ so that $\hat{\mathbf{n}} \times d\hat{\mathbf{n}}$ is essentially $d\theta$ which, on integration, gives a factor 2π . Hence the net angle of Thomas precession during one orbit is given by

$$\Phi = 2\pi(\cosh \chi - 1). \quad (14.8)$$

The similarity between the net angle of turn of the Foucault pendulum and the net Thomas precession angle is now obvious when we compare Eq. (14.8) with Eq. (14.5). We know that in the case of Lorentz transformations, one replaces real angles by imaginary angles which accounts for the difference between the cos and cosh factors. What we need to do is to make this analogy mathematically precise which will be our next task. It will turn out that the sphere and the cone we introduced in the real space, to study the Foucault pendulum, have to be introduced in the velocity space to analyze Thomas precession.

*Velocity space in
relativity*

Before exploring the relativistic velocity space, let us warm-up by asking the following question: Consider two frames S_1 and S_2 which move with velocities \mathbf{v}_1 and \mathbf{v}_2 with respect to a third inertial frame S_0 . What is the magnitude of the relative velocity between the two frames? This is most easily done using Lorentz invariance and four vectors (and to simplify notation we will use units with $c = 1$). We can associate with the 3-velocities \mathbf{v}_1 and \mathbf{v}_2 , the corresponding four velocities, given by $u_1^i = (\gamma_1, \gamma_1 \mathbf{v}_1)$ and $u_2^i = (\gamma_2, \gamma_2 \mathbf{v}_2)$ with all the components being measured in S_0 . On the other hand, with respect to S_1 , this four vector will have the components $u_1^i = (1, 0)$ and $u_2^i = (\gamma, \gamma \mathbf{v})$ where \mathbf{v} (by definition) is the relative velocity between the frames. To determine the magnitude of this quantity, we note that in this frame S_1 we can write $\gamma = -u_{1i}u_2^i$. But since this expression is Lorentz invariant, we can evaluate it in any inertial frame. In S_0 , with $u_1^i = (\gamma_1, \gamma_1 \mathbf{v}_1)$, $u_2^i = (\gamma_2, \gamma_2 \mathbf{v}_2)$ this has the value

$$\gamma = (1 - v^2)^{-1/2} = \gamma_1 \gamma_2 - \gamma_1 \gamma_2 \mathbf{v}_1 \cdot \mathbf{v}_2. \quad (14.9)$$

Simplifying this expression we get

$$v^2 = \frac{(1 - \mathbf{v}_1 \cdot \mathbf{v}_2)^2 - (1 - v_1^2)(1 - v_2^2)}{(1 - \mathbf{v}_1 \cdot \mathbf{v}_2)^2} = \frac{(\mathbf{v}_1 - \mathbf{v}_2)^2 - (\mathbf{v}_1 \times \mathbf{v}_2)^2}{(1 - \mathbf{v}_1 \cdot \mathbf{v}_2)^2}. \quad (14.10)$$

We next consider a 3-dimensional abstract space in which each point represents a velocity of a Lorentz frame measured with respect to some fiducial frame. We are interested in defining the notion of ‘distance’ between two points in this velocity space. Consider two nearby points which correspond to velocities \mathbf{v} and $\mathbf{v} + d\mathbf{v}$ that differ by an infinitesimal quantity. Using the analogy with the usual 3-dimensional flat space, one would have assumed that the “distance” between these two points is just

$$|d\mathbf{v}|^2 = dv_x^2 + dv_y^2 + dv_z^2 = dv^2 + v^2(d\theta^2 + \sin^2 \theta d\phi^2), \quad (14.11)$$

*Relative velocity, by
a simple trick*

*A metric on the
velocity space*

where $v = |\mathbf{v}|$ and (θ, ϕ) denote the direction of \mathbf{v} . In *non-relativistic* physics, this distance also corresponds to the magnitude of the relative velocity between the two frames. However, we have just seen that the relative velocity between two frames in *relativistic* mechanics is different and given by Eq. (14.10). It is more natural to define the distance between two points in the velocity space to be the relative velocity between the respective frames. In that case, the infinitesimal “distance” between the two points in the velocity space will be given by Eq. (14.10) with $\mathbf{v}_1 = \mathbf{v}$ and $\mathbf{v}_2 = \mathbf{v} + d\mathbf{v}$. So

$$dl_v^2 = \frac{(d\mathbf{v})^2 - (\mathbf{v} \times d\mathbf{v})^2}{(1 - v^2)^2} . \quad (14.12)$$

Using the relations

$$(\mathbf{v} \times d\mathbf{v})^2 = v^2(d\mathbf{v})^2 - (\mathbf{v} \cdot d\mathbf{v})^2; \quad (\mathbf{v} \cdot d\mathbf{v})^2 = v^2(dv)^2 , \quad (14.13)$$

and using Eq. (14.11) where θ, ϕ are the polar and azimuthal angles of the direction of \mathbf{v} , we get

$$dl_v^2 = \frac{dv^2}{(1 - v^2)^2} + \frac{v^2}{1 - v^2} (d\theta^2 + \sin^2 \theta d\phi^2) . \quad (14.14)$$

If we use the rapidity χ in place of v through the equation $v = \tanh \chi$, the line element in Eq. (14.14) becomes:

$$dl_v^2 = d\chi^2 + \sinh^2 \chi (d\theta^2 + \sin^2 \theta d\phi^2) . \quad (14.15)$$

This is an example of a *curved* space within the context of special relativity.

If we now change from real angles to the imaginary ones, by writing $\chi = i\eta$, the line element becomes

*This is called the
(three dimensional)
Lobachevsky space*

$$-dl_v^2 = d\eta^2 + \sin^2 \eta (d\theta^2 + \sin^2 \theta d\phi^2) , \quad (14.16)$$

which (except for an overall sign which is irrelevant) represents the distances on a 3-sphere having the three angles η, θ and ϕ as its coordinates.

Of these three angles, θ and ϕ denotes the direction of velocity in the real space as well. When a particle moves in the x-y plane in the real space, its velocity vector lies in the $\theta = \pi/2$ plane and the relevant part of the metric reduces to

$$dL_v^2 = d\eta^2 + \sin^2 \eta d\phi^2 , \quad (14.17)$$

which is just a metric on the 2-sphere. Further, if the particle is moving on a circular orbit having a constant magnitude for the velocity, it follows a curve of $\eta = \text{constant}$ on this 2-sphere. This completes the analogy with the Foucault pendulum, which moves on a constant latitude curve. If the

*Parallel transport,
now in velocity
space*

particle carries a spin, the orbit will transport the spin vector along this circular orbit. As we have seen earlier, the orientation of the vector will not coincide with the original one when the orbit is completed and we expect a difference of $2\pi(1 - \cos \eta) = 2\pi(1 - \cosh \chi)$. So the magnitude of the Thomas precession, over one period is given precisely by Eq. (14.8).

When one moves along a curve in the velocity space, one is sampling different (instantaneously) co-moving Lorentz frames obtained by Lorentz boosts along different directions. As we saw in the last chapter, Lorentz boosts along different directions do not, in general, commute. This leads to the result that if we move along a closed curve in the velocity space (treated as representing different Lorentz boosts) the orientation of the spatial axes would have changed when we complete the loop.

*A nice result in
differential geometry*

The ideas described above are actually of far more general validity. Whenever a vector is transported around a closed curve on the surface of a sphere, the net change in its orientation can be related to the solid angle subtended by the area enclosed by the curve. In the case of the Foucault pendulum, the relevant vector describes the orientation of the plane of the pendulum and the transport is around a circle on the surface of the Earth. In the case of Thomas precession, the relevant vector is the spin of the particle and the transport occurs in the velocity space. Ultimately, both the effects — the Foucault pendulum and Thomas precession — arise because the corresponding space in which the vector is being transported (surface of Earth, relativistic velocity space, respectively) is curved.

Why this Appendix?

Appendix: There is an elegant and geometrical way of obtaining all these results using the concept of parallel transport of vectors on a sphere. Though somewhat more advanced than the other concepts developed here, its sheer elegance makes the case for its inclusion.

Consider moving a vector around a closed curve C “always parallel to itself”. (It is this notion which we will not bother to make precise at this stage, and will rely on your intuition!) If the closed curve was drawn on a sheet of paper and you did this, the vector will point in the original direction when it completes the circuit around the curve. What happens if the closed curve C was drawn on the surface of a sphere? Now the direction of the vector will not coincide with the original direction when it completes the loop. It would have rotated by an angle

$$\alpha(C) = \frac{S(C)}{r^2}, \quad (14.18)$$

The central result

where $S(C)$ is the area enclosed by the curve. There are several ways of proving this result, but probably the most intuitive one is the following:

*1. Approximate the
closed curve by a
polygon with large
number of sides*

We first note that any closed curve can be approximated by a polygon of N sides, with very large N , to as much accuracy as we want. This is clear when the curve is drawn on a sheet of paper when the sides of the polygon is made of usual straight lines. To do the same on the surface of a sphere,

we need the notion of a straight line on the sphere. We know that the curve of shortest distance between any two points on the surface of a sphere is the relevant arc of a great circle, which is the circle passing through the two points with its center at the center of the sphere, and radius equal to the radius of the sphere. [Proof: The shortest distance between two points on the equator is clearly the minor arc along the equator. Given any two points on the surface of a sphere, you can draw an “equator” through them!]. Thus, arcs of great circles generalize the notion of straight lines to the surface of a sphere. So, if we can prove Eq. (14.18) for a trip around a large N -gon on the sphere, with sides made of the arcs of great circles, we are done.

We next note that you can divide up any large polygon into triangles. Again, this fact is obvious if the polygon is on a sheet of paper. To do it on a sphere, we have to generalize the notion of a triangle on to the surface of a sphere. This is easy because the triangular region is bounded by three straight lines and we already know how to define a straight line on the surface of a sphere. It is therefore natural to define a triangle in terms of three intersecting great circles. The area of the polygon is, of course, the sum of the areas of the triangles it is decomposed into. We can now think of moving the vector around the polygon as equivalent to moving it around the individual triangles of which the polygon is made of. Both $\alpha(C)$ and $S(C)$ in Eq. (14.18) add up to give this result. Thus, if we can prove Eq. (14.18) for moving a vector around a triangle drawn on the surface of a sphere, we are done.

2. *Divide the polygon into triangles*

Let us first compute the angle by which the vector rotates when taken around a triangle. Nothing happens to the vector’s orientation when it is moving along the straight lines, being either parallel or perpendicular to the line. All the rotations occur at the three vertices. It is easy to see that, if the three angles of the triangle are $(\theta_1, \theta_2, \theta_3)$ then the rotations are by the amounts $(\theta_1 - \pi, \theta_2 - \pi, \theta_3 - \pi)$ so that the total rotation is by an angle $(\theta_1 + \theta_2 + \theta_3 - 3\pi)$. Since 2π doesn’t count, this is same as a rotation by the angle

3. *Find the result for a triangle and you are done!*

$$\alpha(C) = (\theta_1 + \theta_2 + \theta_3 - \pi) . \quad (14.19)$$

What we need to do is to relate this to the area of the triangle.

This is easy to do. In Fig. 14.4 we take one of the sides of the triangle, AB, and extend it to form the great circle. The “northern hemisphere” formed by this great circle has an area $2\pi r^2$. Similarly, note that the triangle we are interested in (with area S) and the adjacent triangle (S') together form a lune of a sphere. Its area will be a fraction $(\theta_1/2\pi)$ of the full sphere. That is, $2\theta_1 r^2$. Elementary addition of the areas now give us the relation $2\pi r^2 = 2\theta_1 r^2 + 2\theta_2 r^2 + 2\theta_3 r^2 - 2S$. Re-arranging and using Eq. (14.19) we get the required relation

$$S(C) = (\theta_1 + \theta_2 + \theta_3 - \pi)r^2 = \alpha(C)r^2; \quad \alpha(C) = \frac{S(C)}{r^2} . \quad (14.20)$$

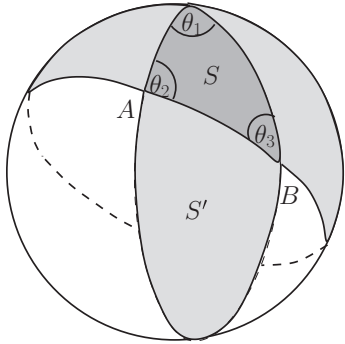


Fig. 14.4: Relating the area $S(C)$ to the angles of the spherical triangle; see text for the discussion.

*Relation to the
Foucault pendulum*

After all this preamble and elegant geometry, let us get back to Foucault and Thomas. The plane of rotation of the Foucault pendulum defines a vector which is normal to the plane. When the pendulum goes around the Earth due to Earth’s rotation, this vector makes a circuit at a fixed latitude. Of course, a given latitude defines a simple curve C on the surface of the sphere, viz., a minor circle with the center located on the axis of rotation. The area $S(C)$ of the sphere enclosed by this curve of constant co-latitude θ is simple to compute and is given by $S(C) = 2\pi r^2 - 2\pi r^2 \cos \theta$. From our result in Eq. (14.18), it follows that the vector defining the plane of the Foucault pendulum will rotate by the amount

$$\alpha(C) = \frac{S(C)}{r^2} = 2\pi - 2\pi \cos \theta \rightarrow -2\pi \cos \theta \; , \tag{14.21}$$

when we ignore the 2π factor. All that the Earth does is to parallel transport the vector defining the normal to the plane of the Foucault pendulum around a circle of constant latitude λ ! (If you use the colatitude θ , the sines become cosines.)

*Do the same in
velocity space and
you get Thomas
precession*

One can do all these in the velocity space as well — which is a pseudo sphere rather than a sphere. A particle moving in a closed orbit in real space will trace a closed curve in the velocity space as well. It is also possible to define the motion of a suitable vector — like the normal to the plane of the Foucault pendulum — in this case too. The Thomas precession is then related to the net rotation of this vector when it is dragged around a closed curve in the velocity space. Because of the similarity between the sphere and the pseudo sphere, a relation similar to Eq. (14.18) holds in this case as well. By calculating the relevant areas using the metric in the velocity space, we can once again obtain the expression for the Thomas precession. You may want to have some fun, filling in the details yourself.

Let us begin by discussing the — apparently elementary — situation of the transition from special relativity to non-relativistic mechanics (NRM) by taking the limit $c^{-1} \rightarrow 0$. (This involves moving along SR to NRM in Fig. 1.1.) While the text books consider this limiting procedure as straightforward, we will see that some curious features arise [2] when we evaluate the limiting form of the action functional in this context.

It is not as simple as you might think

We know that special relativistic mechanics is invariant under a Lorentz transformation of the coordinates, while non-relativistic mechanics is invariant under a Galilean transformation of the coordinates ($x' = x - Vt$, $t' = t$). Given the fact that one recovers the Galilean coordinate transformation by setting $c = \infty$ in the Lorentz transformation equations, one would have thought that any theory which is Lorentz invariant will lead to a theory which is invariant under Galilean transformations in the limit of $c \rightarrow \infty$. As we shall see, it is not so simple.

To illustrate what is involved, let us begin with the One-Body-Problem in physics, viz. the description of a free particle by a suitable action functional in special relativity and compare it with the situation in non-relativistic mechanics. The action \mathcal{A} is, in general, given by

Form of free particle Lagrangian

$$\mathcal{A} = \int L(\mathbf{x}, \mathbf{v}, t) dt . \quad (15.1)$$

But, for a free particle, all locations and directions in space are equivalent; so are all moments of time. If the free particle Lagrangian has been invariant under space and time translations and rotations, it can only be a function of the square of the particle's velocity; i.e., $L = L(v^2)$. *This holds both in the case of relativistic and non-relativistic mechanics.* It is, however, impossible to proceed further and determine the explicit form of $L(v^2)$, without making some additional assumptions. We now have to make a distinction between non-relativistic and relativistic mechanics by postu-

lating the invariance of physical laws under different sets of coordinate transformations.

Symmetry: Galilean invariance

Let us first consider the non-relativistic theory. Here, we postulate that the equations of motion should retain the same form when we make a Galilean transformation:

$$x = x' + Vt; \quad t = t', \quad (15.2)$$

from the co-ordinates (x, t) of an inertial frame S to the co-ordinates (x', t') of another frame S' moving with a uniform velocity V along the positive x -direction with respect to S . (In this and what follows, we suppress the two spatial dimensions and work in $(1 + 1)$ -dimensions for simplicity.) The corresponding velocity transformation is $v = v' + V$ where v and v' are the velocities measured in frames S and S' respectively.

The Lagrangian flunks the test!

The sufficient (though not necessary) condition for the equations of motion to retain the same form in both S and S' is that the action should be invariant under the transformations in Eq. (15.2). It is, however, straightforward to see that *no* non-trivial function $L(v^2)$ has this property! We *cannot* construct an action functional which remains invariant under the Galilean transformation. Rather surprisingly, no Lagrangian exists which respects isotropy and homogeneity of space and strict invariance under Galilean transformations.

Conventional wisdom

Maybe this should warn us that something is wrong and maybe we should have abandoned the Galilean transformation! But historically, one took the easy way out by noting that the equations of motion will remain invariant even if the Lagrangian is not, as long as the Lagrangian changes *only* by the addition of a total time derivative of a function of coordinates and time. The trouble with this option is that, while it is fine in a classical theory, quantum mechanics cares (in the path integral approach) about the exact numerical value of the action. Since classical theories are approximate, and nature is quantum mechanical, we should expect trouble. But let us ignore all this for a moment and proceed further along conventional lines.

It is easy to show that, with these relaxed conditions, we can use a Lagrangian that is *proportional* to the square of the velocity i.e. $L \propto v^2$ or $L = (1/2)mv^2$ where m is defined to be the mass of the particle. In this case, the Lagrangians L and L' , in the two frames of reference S and S' , *differ* by a total time derivative of a function of co-ordinates and time:

$$L = \frac{1}{2}mv^2 = \frac{1}{2}m(v' + V)^2 = L' + \frac{d}{dt}(mx'V + \frac{1}{2}mV^2t). \quad (15.3)$$

The corresponding actions differ by contributions at the end points:¹

$$\mathcal{A} = \mathcal{A}' + \left(mx'V + \frac{1}{2}mV^2t' \right) \Big|_1^2. \quad (15.4)$$

Also note that the canonical momentum ($p' = p - mV$) and the energy ($E' = E - pV + (1/2)mV^2$) are *not* invariant when we transform from L to L' .

Exact theories are more beautiful than approximate ones and show greater level of symmetry. In this context, this is precisely what happens when we proceed from non-relativistic mechanics to relativistic mechanics. In special relativity, we replace Eq. (15.2) by the Lorentz transformations between S and S' of the form:

It is much nicer for the relativistic free particle

$$x = \frac{x' + Vt'}{(1 - V^2/c^2)^{1/2}}; \quad t = \frac{t' + Vx'/c^2}{(1 - V^2/c^2)^{1/2}}. \quad (15.5)$$

The transformation of velocities is now given by:

$$v = \frac{v' + V}{1 + v'V/c^2}. \quad (15.6)$$

It is now possible to construct an action functional which is *actually invariant* (instead of picking up an extra boundary term) under the transformations in Eq. (15.5). This is given by (see Chapter 2):

$$\mathcal{A} = \alpha \int (1 - v^2/c^2)^{1/2} dt, \quad (15.7)$$

where α is a constant.

But it is simply *not* possible to choose α such that Eq. (15.7) reduces to the action for non-relativistic mechanics when $c \rightarrow \infty$! The best we can do is to choose α in such a way that in the non-relativistic limit, we get back the non-relativistic form of the action, *apart from a constant term* in the Lagrangian. This amounts to the standard choice of $\alpha = -mc^2$. Hence, in special relativity, the action for a free particle is taken to be:

$$\mathcal{A} = -mc^2 \int (1 - v^2/c^2)^{1/2} dt. \quad (15.8)$$

The above text book discussion, however, raises some issues.

¹ As an aside, we note the following amusing fact: We implemented homogeneity in time and space in the free particle Lagrangian by excluding the explicit dependence of L on \mathbf{x} or t but incorporated Galilean invariance by allowing L to pick up a total derivative. It is possible to do the converse. One can write down free particle Lagrangians which are *strictly invariant* under Galilean transformations but differ from the standard Lagrangian $L_0 = (1/2)mv^2$ by a total time derivative. A simple example is $L = (1/2)m(\mathbf{v} - \mathbf{x}/t)^2$ which is *invariant* under a Galilean transformation but depends on t and \mathbf{x} . However, L differs from $L_0 = (1/2)mv^2$ by the total time derivative $-d/dt((1/2)mx^2/t)$ which shows that the dependence on t and \mathbf{x} is of no consequence.

The troublesome questions

We expect the non-relativistic theory to arise as a limiting case of the fully relativistic theory, in the limit of $c \rightarrow \infty$. The Lorentz transformation equations *do* reduce (strictly) to the Galilean transformation equations in the limit of $c \rightarrow \infty$. However, the special relativistic action does *not* reduce to the non-relativistic action in this limit, but instead picks up an extra term, $-mc^2t$ evaluated at the end points. This fact is a bit surprising by itself. As we shall see, this term — which is usually ignored in textbooks as being due to “just an addition of a constant to a Lagrangian” — has some interesting implications for the structure of special relativity and non-relativistic mechanics. This is already apparent from the fact that the relativistic action in Eq. (15.8) blows up in the limit of $c \rightarrow \infty$ and does not have a valid limit at all.

Don't even think about it!

You might think one can “renormalize” this action by adding a term $\mathcal{A}_1 \equiv mc^2t$ to Eq. (15.8), then $\mathcal{A} + \mathcal{A}_1$ will have a proper limit. Perish the thought! If you do that, the action will *not* be Lorentz invariant (because \mathcal{A}_1 is not), and hence this “renormalization” is illegal. We will see repeatedly that the term mc^2t plays a crucial role in our future discussion.

The real issue is in quantum theory

Such issues are usually ignored by noting that the equations of motion do not change when a total time derivative of a function (of coordinates and time) is added to the Lagrangian, and hence such action functionals are equivalent as far as physical phenomena are concerned. As I said before, this is true in *classical* physics, but in *quantum* theory, the value of the action is closely related to the phase of the wavefunction (see Chapter 2). Our result shows that the phase of the wavefunction of a free particle remains invariant under Lorentz transformations, but in the $c \rightarrow \infty$ limit, this *invariance gets broken*.

Schrödinger equation in a non-inertial frame

Since the issue really arises only in the quantum theory, we need to examine how the Schrödinger equation transforms under the Galilean transformation. To get more insight (and with future applications in mind) we will work out a slightly more general case of one-body-problem in a *non-inertial* frame. We will consider the transformation from a frame of reference $S = (t, x)$ to a frame $S' = (t, x') \equiv (t, x - \xi(t))$ where $\xi(t)$ is an arbitrary function. When $\xi(t) = Vt$ this corresponds to standard Galilean transformation while for a general $\xi(t)$ it describes a transformation to a non-inertial frame with acceleration $\ddot{\xi}$. The Lagrangian of the free particle in S' is given by

$$L' = \frac{1}{2}m\dot{x}'^2 = \frac{1}{2}m\dot{x}^2 - m\dot{x}\dot{\xi} + \frac{1}{2}m\dot{\xi}^2, \quad (15.9)$$

which can be written in a physically more meaningful way as:

$$L' = L + \frac{df}{dt}, \quad (15.10)$$

and L is a new Lagrangian given by

$$L = \frac{1}{2}m\dot{x}^2 + m\ddot{\xi}x, \quad (15.11)$$

and

$$f \equiv -m\ddot{\xi}x + \int \frac{1}{2}m\ddot{\xi}^2 dt. \quad (15.12)$$

Clearly, L is equivalent to L' as far as the equations of motion are concerned, since the total time derivative df/dt does not contribute to the equations of motion. Moreover, L represents the Lagrangian for a particle acted upon by a force $m\ddot{\xi}$ or, equivalently, a particle located in a spatially homogeneous (but time dependent) gravitational field $\ddot{\xi}$. This is precisely what we would have expected from the Principle of Equivalence; so everything makes sense.

By and large, it goes as expected

But when we add $df(x,t)/dt$, to the Lagrangian L , thus transforming it to $L' = L + df/dt$, both the canonical momentum and the Hamiltonian change, becoming

$$p' = p + \frac{\partial f}{\partial x}; \quad H' = H - \frac{\partial f}{\partial t}. \quad (15.13)$$

In quantum mechanics, the time evolution of the wavefunction is determined by the Hamiltonian operator and hence, the form of the wavefunction must change when we make a co-ordinate transformation from S to S' . Let $\Psi'(t, x')$ be the quantum-mechanical wavefunction for the free particle in the frame S' . Then, it can be shown that the corresponding wavefunction $\Psi(t, x)$ for the same particle in the frame S is given by:

The ψ is not a scalar under the coordinate transformation!

$$\Psi(t, x) = \Psi'(t, x - \xi(t))e^{-if/\hbar}, \quad (15.14)$$

where $f(t, x)$ is given by Eq. (15.12) and $\Psi(t, x)$ satisfies the equation

$$i\hbar \frac{\partial \Psi(t, x)}{\partial t} = -\frac{\hbar^2}{2m} \frac{\partial^2 \Psi(t, x)}{\partial x^2} - m\ddot{\xi}x\Psi(t, x), \quad (15.15)$$

in the frame of reference S . (Equation 15.15 is derived in Appendix 1.) We see from Eq. (15.11) that in this frame, the particles do experience a pseudo-force $-m\ddot{\xi}$ which arises from the “pseudo-potential” energy term, $-m\ddot{\xi}x$, and Eq. (15.15) is indeed the Schrödinger equation with a “pseudo-potential” energy term, $-m\ddot{\xi}x$. But Eq. (15.14) tells us that you *cannot* just obtain the wavefunction in this frame by substituting $x' = x - \xi(t)$ which is what we would have done if the wavefunction is a scalar; you *have to change the phase as well by f* . But from Eq. (15.10) we know that the two actions corresponding to L and L' differ by f ; so the phase change is exactly the change in the action.

Box 15.1: Quantum particle in constant gravitational field

The result in Eq. (15.15) might look like one of those formal things but it has practical applications. Note that it allows you to solve a *time dependent* Schrödinger equation in a class of potential of the form $F(t)x$ by just finding the free particle solution and transforming to a different frame!

*Particle in a
constant force field*

As a simple application of this, consider a case of the particle located in a uniform force field with the Hamiltonian $H = (1/2)p^2 + ax$. The usual way of determining the eigenfunctions $H\phi_E = E\phi_E$ leads to Airy functions in x -space. This, however, is one problem in which the momentum space representation of the operators with $x = i(\partial/\partial p)$ turns out to be easier to handle!. The Schrödinger equation in the p -representation is now $ia(\partial\phi/\partial p) = [E - (p^2/2)]\phi$. Integrating this equation and then Fourier transforming we get the solution in the x -representation to be

$$\phi_E(x) = \int_{-\infty}^{\infty} dp \exp i[p(x - E/a) + (p^3/6a)] , \quad (15.16)$$

which is indeed an integral representation for the Airy function (see e.g., Ref. [64]).

*Cleverer way to get
the same result*

But we can solve this problem using Eq. (15.15)! We begin with the simplest free particle solution to the Schrödinger equation, which are the momentum eigenfunctions $\psi_{\text{free}}(t, x) = \exp(-ipx + ip^2t/2)$. We next obtain the solution to Eq. (15.15) by the simple transformation $x \rightarrow x + \ell(t)$ where $\ddot{\ell} = a = \text{constant}$ and the addition of a phase as indicated in Eq. (15.14). This gives the solution:

$$\psi = \exp -i[x(p - at) + (1/2)pat^2 - (1/2)p^2t - (1/6)a^2t^3] . \quad (15.17)$$

This is, of course, not an energy eigenfunction. However, a Fourier transform of this expression with respect to t

$$\phi_E(x) = \int_{-\infty}^{\infty} dt \psi(t, x) \exp iEt , \quad (15.18)$$

will give the energy eigenfunctions for a particle moving in a uniform force field. Changing the variable of integration from t to $\xi \equiv (at - p)$, you will find that various terms cancel out nicely, leading to

$$\phi_E(x) \propto \int_{-\infty}^{\infty} d\xi \exp i[\xi(x - E/a) + (\xi^3/6a)] , \quad (15.19)$$

which are the same energy eigenfunctions as in Eq. (15.16) except for an unimportant phase!

After having obtained the result for a general $\xi(t)$, let us get back to the Galilean transformation, which corresponds to $\xi(t) = Vt$ and

$$f = -mxV + \frac{1}{2}mV^2t. \quad (15.20)$$

So, in this case when we need to relate the two wavefunctions using Eq. (15.14) we get:

$$\Psi(t, x) = \Psi'(t, x - Vt) \exp[(-i/\hbar)(-mxV + (1/2)mV^2t)]. \quad (15.21)$$

That is, we need to transform the wavefunction, treating it as a scalar, *and then add an extra phase* which is consistent with what we found earlier. As we have said before, all this is perfectly consistent as regards the application of the Galilean transformation in quantum mechanics.

The solutions of the Schrödinger equation are not scalars under Galilean transformations ...

Classically we saw that the action was invariant in special relativity, while it picked up an end-point contribution in the non-relativistic case. What is the analogue for the relativistic case when we treat the particle quantum mechanically? In this case, one could use the Klein-Gordon equation to describe a spin-zero particle. Since Klein-Gordon equation is fully Lorentz invariant, its solution will transform as a scalar when we go from one frame to another. No additional phase should appear. If so, how is it that the Klein-Gordon equation is invariant under the Lorentz transformation, but the Schrödinger equation — which is presumably obtained in the $c \rightarrow \infty$ limit of the Klein-Gordon equation — is *not* invariant under the Galilean transformation, given the fact that the Lorentz transformation reduces to the Galilean transformation in the appropriate limit?

... but the solutions of the KG equation are scalars under Lorentz transformations; How come?

This has to do with the manner in which one obtains the Schrödinger equation from the Klein-Gordon equation and brings to the center stage the role of the mc^2 term in the phase. We will outline how the extra phase in Eq. (15.21) can be obtained from a fully invariant Klein-Gordon equation.

Consider the wavefunction $\Phi(t, x)$ which is the solution to a free particle Klein-Gordon equation. We know that under a Lorentz transformation, $\Phi(t, x) \Rightarrow \Phi'(t', x')$, thus transforming as a scalar. To obtain the Schrödinger equation for a wavefunction $\psi(t, x)$ we first have to separate the mc^2t term from the phase of the Φ by writing

$$\Phi(t, x) = \psi(t, x) \exp[-imc^2t]. \quad (15.22)$$

It is then straightforward to show that in the limit of $c \rightarrow \infty$, $\psi(t, x)$ will satisfy a free particle Schrödinger equation. (We will demonstrate a more general result in the presence of a gravitational field later on, and hence we skip the algebraic details here; see Eq. (15.28).) To obtain the Schrödinger equation in S' , we have to similarly write $\Phi'(t', x') = \psi'(t', x') \exp(-imc^2t')$. The fact that Φ transforms as a scalar can now be

The crucial phase difference

used to relate ψ and ψ' , and we find a remarkable result:

$$\psi' = \psi \exp[-imc^2(t-t')]. \quad (15.23)$$

We see that, in addition to the scalar transformation, the wavefunction picks up a phase which is just $mc^2(t-t')$. *Incredibly enough, this expression has a finite, non-zero limit when $c \rightarrow \infty$!* Evaluating this quantity in the limit of $c \rightarrow \infty$, we get

$$\begin{aligned} mc^2(t-t') &= mc^2 \gamma \left(t' + \frac{Vx'}{c^2} \right) - mc^2 t' \\ &= mc^2 \left[t' + \frac{Vx'}{c^2} + \frac{1}{2} \frac{V^2 t'^2}{c^2} + \mathcal{O} \left(\frac{1}{c^4} \right) \right] - mc^2 t' \\ &= mVx' + \frac{mV^2 t'}{2} + \mathcal{O} \left(\frac{1}{c^2} \right). \end{aligned} \quad (15.24)$$

This is precisely the mysterious phase which occurs in the Schrödinger equation under a Galilean transformation! It has a simple interpretation as being equal to $mc^2(t-t')$, thus emphasizing the role of rest energy *even in the non-relativistic limit*. This result tells you the innocuous phase we needed to add to the wavefunction in the case of non-relativistic quantum mechanics actually arises from special relativity and has an elementary interpretation in special relativity. *Once again, more exact theories make better sense than approximate ones!*

You need special relativity to understand non-relativistic physics!

One might think this is probably just a coincidence, but it is not. To see that, let us consider a more complicated situation — not that of uniform motion but the one with an acceleration. We are now looking at the Klein-Gordon and Schrödinger equations for a free particle in non-inertial frames and we want to know whether the phase acquired in non-relativistic quantum mechanics is actually related to the time difference as measured by different clocks. As we will see, this is indeed the case!

The idea works in more general cases

In the relativistic case, we have a quantum scalar field satisfying the free-particle Klein-Gordon equation in one frame of reference (S), which we call (x, t) ; the frame (x, t) being arbitrarily accelerated with time-dependent acceleration $g(t)$ with respect to an inertial coordinate system $S' = (X, T)$. In S (known as the *generalized Rindler frame*) the metric is given by (see the Appendix 2):

Spacetime metric in an accelerated frame

$$ds^2 = - \left(1 + \frac{g(t)x}{c^2} \right) dt^2 + dx^2. \quad (15.25)$$

The explicit co-ordinate transformation (see e.g., [65]) between S and the inertial frame S' is given by Eq. (15.43) in Appendix 2. The Klein-Gordon equation for a scalar field $\Phi(x, t)$ in an arbitrary frame is given by

$$\frac{1}{\sqrt{-g}} \partial_i (\sqrt{-g} g^{ik} \partial_k \Phi) = \mu^2 \Phi; \quad \mu \equiv \frac{mc}{\hbar}. \quad (15.26)$$

(The complicated looking expression is just the \square in curvilinear coordinates. You know that while $\nabla^2 = \partial^2/\partial x^2 + \partial^2/\partial y^2 + \partial^2/\partial z^2$ in Cartesian coordinates, it becomes more complicated in the spherical polar coordinates. What we have here is just a similar result for \square .) Using the form of the metric as given in Eq. (15.25), we can expand this as:

It's really quite simple

$$-\frac{1}{(1+g(t)x)^2} \frac{\partial^2 \Phi}{\partial t^2} + x \frac{dg}{dt} \frac{\partial \Phi}{\partial t} \frac{1}{(1+g(t)x)^3} + \frac{\partial^2 \Phi}{\partial x^2} + \frac{\partial \Phi}{\partial x} \frac{g}{(1+g(t)x)} = \mu^2 \Phi. \quad (15.27)$$

We now substitute $\Phi(x, t) = \psi(x, t)e^{-i\mu t}$ into Eq. (15.27), we get on retaining terms to the lowest order (upto, but excluding, order gx/c^2), the equation:

$$i\hbar \frac{\partial \psi}{\partial t} = -\frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial x^2} + mg(t)x\psi, \quad (15.28)$$

which is identical to the Schrödinger equation for a particle of mass m in an accelerated frame of reference moving with acceleration $-g(t)$ or equivalently, in a time-dependent gravitational field of strength $g(t)$. Hence, we see that the Klein-Gordon equation does reduce, in the appropriate limit, to the Schrödinger equation, with the term $mg(t)x$ indicating the accelerated nature of the frame.

All this is fine, but what about the phase factor? The solution to the Klein-Gordon equation is invariant when we go from S to S' ; i.e., $\Phi'(T, X) = \Phi(t, x)$. But the solution to the Schrödinger equation, Eq. (15.28) acquires an extra phase f given in Eq. (15.14). Where does this come from?

In fact, it has a direct physical meaning. We can transform the free particle solution to the Klein-Gordon equation in the inertial frame, $\Phi(T, X)$, as a scalar to the non-inertial frame, thus obtaining $\Phi(t, x)$. But the *non-relativistic limits* of $\Phi(T, X)$ and $\Phi(t, x)$ will differ by a phase term $mc^2(t - T)$, which, in the appropriate limit, will give the correct phase dependence arrived at in Eq. (15.14) when we consider the *effect of gravitational time dilation!*

Clocks run differently in different frames ...

In the presence of a gravitational potential ϕ , the proper time lapse dT of a co-moving clock is related to the coordinate time lapse dt by (see Chapter 11):

... or in the presence of gravity

$$\begin{aligned} ds^2 &= -c^2 dT^2 = -c^2 \left(1 + \frac{2\phi}{c^2} \right) dt^2 + dx^2 \\ &= -c^2 dt^2 \left[\left(1 + \frac{2\phi}{c^2} \right) - \frac{V^2}{c^2} \right], \end{aligned} \quad (15.29)$$

so that, when $V = \dot{\xi}$, $\phi = x\ddot{\xi}$, we get

$$\begin{aligned} mc^2(t-T) &= -mc^2 \left[\int dt \left(1 - \frac{\dot{\xi}^2}{c^2} + \frac{2x\ddot{\xi}}{c^2} \right)^{1/2} - t \right] \\ &\approx -m \int dt \left(-\frac{\dot{\xi}^2}{2} + x\ddot{\xi} \right) = -mx\dot{\xi} + \frac{1}{2}m \int dt \dot{\xi}^2, \quad (15.30) \end{aligned}$$

which is precisely the phase f found in Eq. (15.12)! Once again, we see that a result in non-relativistic quantum mechanics acquires a simple interpretation when we treat it as *a limit of relativistic theory*, thanks to the factor $mc^2(t-T)$ in the phase. The result also shows that in the instantaneous rest frame of the particle, the phase of the wavefunction evolves as $mc^2 d\tau$, where τ is the proper time shown by the co-moving clock, thereby again validating the principle of equivalence in quantum mechanics.

More exact theories are more elegant and make better sense!

Box 15.2: Why does the harmonic oscillator have coherent states?

You would have learnt that the standard harmonic oscillator admits coherent state solutions in which the probability distribution varies as

$$|\phi_A(t, x)|^2 \propto \exp[-\omega(x - A \cos \omega t)^2]. \quad (15.31)$$

This is obtained by just shifting the ground state probability distribution by $x \rightarrow x - A \cos \omega t$. What is more surprising (in case you did not know) is that such coherent states exist *even for the excited states* of the oscillator with the same shift! (People have tried to find coherent states for other potentials but none of them look as neat as those for the oscillators.) Why does the harmonic oscillator admit such a nice set of states? The existence of such states is a mystery in the conventional approach to quantum mechanics, but our approach based on Eq. (15.9) provides a valuable insight.

To understand this, let us apply the transformation $x \rightarrow \bar{x} = x + \ell(t)$ to the harmonic oscillator Lagrangian $L = (1/2)(\dot{x}^2 - \omega^2 x^2)$. Elementary algebra shows that the new Lagrangian has the structure

$$\bar{L} = (1/2)(\dot{x}^2 - \omega^2 x^2) - (\ddot{\ell} + \omega^2 \ell)x + \frac{df}{dt}, \quad (15.32)$$

where f is again a function determined by $\ell(t)$ but its explicit form is not important. Let us now choose $\ell(t)$ to be a solution to the classical equation of motion $\ddot{\ell} + \omega^2 \ell = 0$. To be specific, we will take $\ell = -A \cos \omega t$. If you want, you can think of this as shifting to a frame

An unexpected bonus

which is oscillating with the particle. We then see that the second term in Eq.(15.21) vanishes and \bar{L} has the *same* form as the original harmonic oscillator Lagrangian except for the total derivative.

The solutions to the Schrödinger equation are, therefore, the *same* as the standard solutions to the harmonic oscillator problem with a shift $x \rightarrow x + \ell(t)$ and an extra phase factor! The probabilities do not care for the phase factor, and we have the result $|\bar{\psi}|^2 = |\psi(x + \ell(t), t)|^2$. If ψ is the ground state, then this shift leads to the standard coherent state. But if you take the n th excited state of the oscillator $\psi_n(x, t)$, shift the coordinate and add a phase, then we get another valid solution $e^{if} \psi_n(x - A \cos \omega t, t)$. As far as the probability goes, $|\psi_n(x - A \cos \omega t, t)|^2$ merely traces the original probability distribution with the mean value oscillating along the classical solution. In our approach, we see that a harmonic oscillator gets mapped *back* to a harmonic oscillator when we move to a frame with $\tilde{\ell} + \omega^2 \ell = 0$ with just a shift in x (and a phase which is irrelevant for the probabilities).

This miracle occurs only for the quadratic potential!

That is why coherent states exist even for the excited states of the harmonic oscillator.

Appendix 1: In this appendix we prove that the wavefunction Eq. (15.14) satisfies the Schrödinger equation, Eq. (15.15). We set $m = \hbar = 1$ for convenience, so that Eq. (15.15) becomes:

$$i \frac{\partial \Psi}{\partial t} = -\frac{1}{2} \frac{\partial^2 \Psi}{\partial x^2} - \ddot{\xi} x \Psi. \quad (15.33)$$

The co-ordinate transformation is given by $x' = x - \xi(t)$, $t' = t$. We now substitute $\Psi(x, t) = \Psi'(x', t') e^{-if}$ into the above equation, where $f = -x\ddot{\xi} + \frac{1}{2} \int \ddot{\xi}^2 dt$. We have the following relations:

$$i \frac{\partial (\Psi' e^{-if})}{\partial t} = i \frac{\partial \Psi'}{\partial t} e^{-if} + e^{-if} \Psi' \frac{\partial f}{\partial t} \quad (15.34)$$

and

$$\frac{\partial}{\partial x} (\Psi' e^{-if}) = e^{-if} \frac{\partial \Psi'}{\partial x'} - i e^{-if} \frac{\partial f}{\partial x} \Psi'. \quad (15.35)$$

Hence,

$$\frac{\partial^2}{\partial x^2} (\Psi' e^{-if}) = e^{-if} \frac{\partial^2 \Psi'}{\partial x'^2} - 2i \frac{\partial \Psi'}{\partial x'} \frac{\partial f}{\partial x} e^{-if} - e^{-if} \left(\frac{\partial f}{\partial x} \right)^2 \Psi', \quad (15.36)$$

where we have used the facts that $\partial \Psi' / \partial x = \partial \Psi' / \partial x'$ and $\partial^2 f / \partial x^2 = 0$. Using these relations, Eq. (15.33) becomes:

$$i \frac{\partial \Psi'}{\partial t} + \Psi' \frac{\partial f}{\partial t} = -\frac{1}{2} \frac{\partial^2 \Psi'}{\partial x'^2} + i \frac{\partial \Psi'}{\partial x'} \frac{\partial f}{\partial x} + \frac{1}{2} \left(\frac{\partial f}{\partial x} \right)^2 \Psi' - \ddot{\xi} \Psi' x. \quad (15.37)$$

We also know that

$$\frac{\partial \Psi'}{\partial t} = \frac{\partial \Psi'}{\partial t'} - \xi \frac{\partial \Psi'}{\partial x'} \quad (15.38)$$

and

$$\frac{\partial f}{\partial x} = -\xi; \quad \frac{\partial f}{\partial t} = -\xi x + \frac{1}{2} \xi^2. \quad (15.39)$$

Using the above relations in Eq. (15.37), it readily transforms to:

$$i \frac{\partial \Psi'}{\partial t'} = -\frac{1}{2} \frac{\partial^2 \Psi'}{\partial x'^2}, \quad (15.40)$$

which is satisfied identically, since we know that $\Psi'(x', t')$ is a solution to the free particle Schrödinger equation in the (x', t') frame of reference. Hence, we see that the wavefunction in Eq. (15.14) satisfies Eq. (15.15).

Appendix 2: We will indicate how to obtain the coordinate system and the metric for an observer moving with an arbitrary, *time dependent* acceleration along the x -axis.

Consider an accelerated observer with the trajectory $T = h(\tau)$, $X = f(\tau)$ and a coordinate velocity $u(\tau) \equiv df/dh$ where τ is the proper time. At any given instant, there exists a Lorentz frame (t, \mathbf{x}) with: (a) the three coordinate axes coinciding with the axes of the accelerating observer, and (b) the origin coinciding with the location of the observer. The Lorentz transformations (with suitable translation of origin) from the global inertial frame coordinates (T, X) to this instantaneously comoving frame is given by (with $c = 1$)

$$X - f(\tau) = \gamma(u)(x + ut); \quad T - h(\tau) = \gamma(u)(t + ux). \quad (15.41)$$

We now define the coordinates for the accelerated observer such that, at $t = 0$ the coordinate labels in the accelerated frame coincide with those in the comoving Lorentz frame. This gives

$$X = f(\tau) + \gamma(u)x; \quad T = h(\tau) + \gamma(u)ux. \quad (15.42)$$

This result can be rewritten in a more explicit form as:

$$\begin{aligned} X &= \int' \sinh \chi(t) dt + x \cosh \chi(t) = \int dt [1 + g(t)x] \sinh \chi(t) \\ T &= \int' \cosh \chi(t) dt + x \sinh \chi(t) = \int dt [1 + g(t)x] \cosh \chi(t), \end{aligned} \quad (15.43)$$

where the function $\chi(t)$ is related to the time dependent acceleration $g(t)$ by $g(t) = (d\chi/dt)$.

We can now find the corresponding metric in the accelerated frame by computing $-dX^2 + dT^2$ in terms of dx and dt . This calculation shows that

the line element in these coordinates is remarkably simple and is given by

$$ds^2 = -(1 + g(t)x)^2 dt^2 + (dx^2 + dy^2 + dz^2) . \quad (15.44)$$

It is amazing that such a simple expression can be obtained for an arbitrary acceleration $g(t)$. When the acceleration is constant, it reduces to the expressions used in the main text.

The Straight and Narrow Path of Waves

16

The unification of electricity and magnetism through Maxwell's equations led to our understanding of light as an electromagnetic wave. This historical milestone allowed us to think of light as made of oscillating electric and magnetic fields, each of which obeys a wave equation. In this chapter we want to look at the wave nature of light from a particular point of view [66] which we will connect up with a seemingly different phenomenon in Chapter 17.

Warm up for the next chapter

For our purpose the vector nature of the electromagnetic field is not relevant (since we will not be interested, e.g., in the polarization of the light). Hence, we will just deal with one component — called $A(t, \mathbf{x})$, say, — of the relevant vector field which satisfies the wave equation. The solution to the wave equation $\square A = 0$ is described by the (real and imaginary parts of the) function $\exp i[\mathbf{k} \cdot \mathbf{x} - \omega t]$. Here \mathbf{k} denotes the direction of propagation of the wave which also determines its frequency through the dispersion relation $\omega = |\mathbf{k}|c$. Since the wave equation is linear in A , superposition of the solutions with different values of \mathbf{k} , each with an amplitude $F_1(\mathbf{k})$, say, leads to:

Scalars will do, nicely

$$A(t, \mathbf{x}) = \int F_1(\mathbf{k}) e^{i\mathbf{k} \cdot \mathbf{x}} e^{-i\omega t} \frac{d^3 k}{(2\pi)^3} . \quad (16.1)$$

We now specialize to a situation which arises in the study of optical phenomenon. Quite often, we are concerned with waves which are propagating broadly along some given direction, say, along the positive z -axis. For example, consider the study of diffraction by a circular hole in a screen which is located in the $z = 0$ plane. We will consider, in such a context, light incident on the screen from the left and getting diffracted; the propagation being essentially along the z -axis with a diffraction spread in the transverse direction. Mathematically, this means that the function $F_1(\mathbf{k})$ is nonzero only for wave vectors with $k_z > 0$. Further, since the wave has a definite frequency ω , the magnitude of the wave vector is fixed at the value ω/c . It follows that one of the components of the wave vector, say

Practical case

k_z , can be expressed in terms of the other three. So, the function F_1 has the structure:

$$F_1(k_z, \mathbf{k}_\perp) = 2\pi f(\mathbf{k}_\perp) \delta_D \left(k_z - \sqrt{\omega^2/c^2 - \mathbf{k}_\perp^2} \right), \quad (16.2)$$

where the subscript \perp denotes the components of the vector in the transverse $x-y$ plane and $f(\mathbf{k}_\perp)$ is an arbitrary function of \mathbf{k}_\perp . Note that — in general — we could have had

$$k_z = \pm \sqrt{\omega^2/c^2 - \mathbf{k}_\perp^2}, \quad (16.3)$$

and we have consciously picked out one with $k_z > 0$ leading to propagation in the direction of positive z -axis.

Substituting this expression in Eq. (16.1), we find that $A(t; z, \mathbf{x}_\perp)$ can be written in the form $a(z, \mathbf{x}_\perp) e^{-i\omega t}$ (in which the oscillations in time have been separated out) where

$$a(z, \mathbf{x}_\perp) = \int \frac{d^2 \mathbf{k}_\perp}{(2\pi)^2} f(\mathbf{k}_\perp) e^{i\mathbf{k}_\perp \cdot \mathbf{x}_\perp} \exp \left[\frac{iz}{c} \sqrt{\omega^2 - c^2 k_\perp^2} \right]. \quad (16.4)$$

Since the time variation of a monochromatic wave is always $\exp(-i\omega t)$, we shall ignore this factor and concentrate on the spatial dependence of the amplitude, $a(z, \mathbf{x}_\perp)$.

*The context of
paraxial optics*

To proceed further, we consider the case in which all the components building up the wave are traveling essentially along the positive z -axis with a small transverse spread. For such a wave traveling, by and large, along the z direction, the transverse components of \mathbf{k} are small compared to its magnitude; that is, $c^2 k_\perp^2 \ll \omega^2$. Using the Taylor series

$$\sqrt{\omega^2 - c^2 k_\perp^2} \cong \omega \left(1 - \frac{1}{2} \frac{c^2 k_\perp^2}{\omega^2} \right) = \omega - \frac{1}{2} \frac{c^2 k_\perp^2}{\omega}, \quad (16.5)$$

in Eq. (16.4), we get:

$$a(z, \mathbf{x}_\perp) \cong e^{i\omega z/c} \int \frac{d^2 \mathbf{k}_\perp}{(2\pi)^2} f(\mathbf{k}_\perp) \exp \left[i(\mathbf{k}_\perp \cdot \mathbf{x}_\perp - (c/2\omega) k_\perp^2 z) \right]. \quad (16.6)$$

This equation describes the propagation of a wave along the positive z -axis with a small spread in the transverse direction. The function $f(\mathbf{k}_\perp)$ can be determined by a simple Fourier transform if the amplitude $a(z', \mathbf{x}'_\perp)$ at some location z' is known. Doing this, we can relate the amplitudes of the wave at two planes with coordinates z and z' by

$$a(z, \mathbf{x}_\perp) = e^{i\omega(z-z')/c} \int d^2 \mathbf{x}'_\perp a(z', \mathbf{x}'_\perp) G(z-z'; \mathbf{x}_\perp - \mathbf{x}'_\perp), \quad (16.7)$$

where

$$G(z - z'; \mathbf{x}_\perp - \mathbf{x}'_\perp) = \int \frac{d^2 \mathbf{k}_\perp}{(2\pi)^2} e^{i\mathbf{k}_\perp \cdot (\mathbf{x}_\perp - \mathbf{x}'_\perp)} e^{-(ic/2\omega)k_\perp^2(z - z')} \\ = \left(\frac{\omega}{2\pi ic} \right) \frac{1}{|z - z'|} \exp \left[\frac{i\omega}{2c} \frac{(\mathbf{x}_\perp - \mathbf{x}'_\perp)^2}{(z - z')} \right] \quad (16.8)$$

The function G may be thought of as a propagator which propagates the amplitude from the location (z', \mathbf{x}'_\perp) to the location (z, \mathbf{x}_\perp) . The factor $e^{i\omega(z - z')/c}$ in Eq. (16.7) does not contribute to the intensity (which is proportional to $|a(z, \mathbf{x}_\perp)|^2$) and we will drop it when not necessary.

Takes it from there to here

Some thought shows that we have achieved something quite extraordinary. We know that the wave amplitude satisfies a *second* order differential equation (viz. the wave equation) and hence its evolution cannot be determined by just knowing the amplitude (i.e., one single function, $a(z', \mathbf{x}'_\perp)$) at a given plane (z', \mathbf{x}'_\perp) . But that is exactly what we have done! This was possible in Eq. (16.7) because of the assumption that the wave is traveling only forward in the z direction. The actual form of the propagator depended on the assumption that the transverse components of the wave vector were small compared to k_z . The study of wave propagation under these approximations is called *paraxial optics*.

This is quite nontrivial

Let us take a closer look at the structure of the propagator G in Eq. (16.8) which introduces a factor $|z - z'|^{-1}$ to the amplitude and — more importantly — contributes an amount

$$\phi = \frac{\omega}{2c} \frac{(\mathbf{x}_\perp - \mathbf{x}'_\perp)^2}{(z - z')} \quad (16.9)$$

to the phase. The change in the amplitude merely reflects the r^{-2} fall off of the intensity (which is proportional to the square of the amplitude) of the wave. But what is the meaning of the phase factor? To understand the origin of the change in phase, note that a path difference Δs between two points in space will introduce a phase difference of $k\Delta s$ in a propagating wave. In our case, it is clear that the phase difference is

Amplitude is easy

$$k\Delta s = \frac{\omega}{c} \left[\sqrt{(\mathbf{x}_\perp - \mathbf{x}'_\perp)^2 + (z - z')^2} - (z - z') \right] \cong \frac{\omega}{c} \left[\frac{1}{2} \frac{(\mathbf{x}_\perp - \mathbf{x}'_\perp)^2}{(z - z')} \right] \quad (16.10)$$

provided the transverse displacements are small compared to the longitudinal distance — an assumption which is central to paraxial optics. With hindsight, we could have guessed the form of G without doing any algebra! In paraxial optics, it introduces a phase corresponding to the path difference and decreases the amplitude to take care of the normal spread of the wave.

Phase difference is path difference

*What optical
systems do*

Equation (16.7) allows us to compute the wave amplitude at any location on the plane $z = z_2$, if the amplitude on a plane $z = z_1 < z_2$ is given. As an application, we now consider a standard situation which arises quite often in optics. A wave front propagates freely up to a plane $z = z_1$ where it passes through an optical system (say a lens, screen with a hole, atmosphere, etc...) which modifies the wave in a particular fashion. The optical system extends from $z = z_1$ to $z = z_2$ and the wave propagates freely for $z > z_2$. We will be interested in the amplitude at $z > z_2$, given the amplitude at $z < z_1$.

It is clear that our Eq. (16.7) can be used to propagate the amplitude from some initial plane $z = z_O < z_1$ to $z = z_1$ and from $z = z_2$ to some final plane $z = z_I > z_2$. (The subscripts *O* and *I* stand for the object and the image, based on the idea of the optical system being a lens.) The propagation of the wave from z_1 to z_2 depends entirely on the optical system and — in fact — defines the particular optical system. An optical system is called *linear* if the output is linear in the input. In such a case, the amplitude at the exit point of the optical system is related to the amplitude at the entrance point by a relation of the form:

$$a(z_2, \mathbf{x}_2) = \int d^2 \mathbf{x}_1 P(z_2, z_1; \mathbf{x}_2, \mathbf{x}_1) a(z_1, \mathbf{x}_1) , \quad (16.11)$$

where the functional form of P determines the kind of optical system. (Here, and in what follows, we shall omit the subscript \perp with the understanding that the vector \mathbf{x} is in the transverse plane and is two dimensional.) In this case, the amplitude at the image plane can be expressed in terms of the amplitude at the object plane by the relation

Complete solution

$$a(z_I, \mathbf{x}_I) = \int d^2 \mathbf{x}_O \mathcal{G}(z_I, z_O; \mathbf{x}_I, \mathbf{x}_O) a(z_O, \mathbf{x}_O) , \quad (16.12)$$

where

$$\begin{aligned} \mathcal{G}(z_I, z_O; \mathbf{x}_I, \mathbf{x}_O) = & \int d^2 \mathbf{x}_2 d^2 \mathbf{x}_1 G(z_I - z_2, \mathbf{x}_I - \mathbf{x}_2) P(z_2, z_1; \mathbf{x}_2, \mathbf{x}_1) \\ & \times G(z_1 - z_O, \mathbf{x}_1 - \mathbf{x}_O) . \end{aligned} \quad (16.13)$$

Given the properties of any linear optical system, one can compute the quantity P , and thus evaluate \mathcal{G} and determine the properties of wave propagation.

*Example: A convex
Lens*

As a simple example, let us compute the form of the function P for a convex lens. If the lens is sufficiently thin, P will be nonzero only at the plane of the lens $z_2 = z_1 = z_L$. Since the lens does not absorb radiation, it cannot change the amplitude $|a(z_L, \mathbf{x}_L)|$ of the incident wave and can only modify the phase. Therefore, P must have the form $P = \exp[i\theta(\mathbf{x}_L)]$. Then

the amplitude at the image plane is given by:

$$\begin{aligned} a(z_I, \mathbf{x}_I) &= \int d^2\mathbf{x}_L a(z_L, \mathbf{x}_L) P(z_L, \mathbf{x}_L) G(z_I - z_L, \mathbf{x}_I - \mathbf{x}_L) \\ &= a \int d^2\mathbf{x}_L e^{i\theta(z_L, \mathbf{x}_L)} G(z_I - z_L, \mathbf{x}_I - \mathbf{x}_L), \end{aligned} \quad (16.14)$$

where we have used the fact that the amplitude $a(z_L, \mathbf{x}_L)$ on the lens plane is constant for a plane wave incident from a large distance. To determine the form of $\theta(\mathbf{x}_L)$, we use the basic defining property of lens of focal length f : If a plane wavefront of constant intensity is incident on the lens plane $z = z_L$, the rays will be focused at a point $z_I = z_L + f$, when the wave nature of the light is ignored. In the limit of zero wavelength of the wave, most of the contributions to the integral come from points at which the phase of the integrand in Eq. (16.14) is stationary. Since the phase of G is $(k/2)[(\Delta\mathbf{x})^2/\Delta z]$, the principle of stationary phase gives the equation,

Plane to sphere

$$\frac{\partial \theta}{\partial \mathbf{x}_L} = \frac{k}{f} (\mathbf{x}_I - \mathbf{x}_L), \quad (16.15)$$

where $f = z_I - z_L$. For the image to be formed along the z -axis, this equation should be satisfied for $\mathbf{x}_I = 0$. Setting $\mathbf{x}_I = 0$, and integrating this equation, we find that $\theta = (-kx_L^2/2f)$ and

$$P(\mathbf{x}_L) = \exp\left(-\frac{ik}{2f}x_L^2\right). \quad (16.16)$$

Thus the effect of a lens is to introduce a phase variation which is quadratic in the transverse coordinates. Such a lens will focus the light to a point on the z -axis, in the limit of zero wavelength.

The most sophisticated way of defining a lens

A geometrical interpretation of this result is given in Fig. 16.1. The constant phase surfaces are planes to the left of the lens, and are arcs of circles (centered on the focus F) to the right of the lens. Changing the constant phase surfaces from the plane to a circle (of radius f) through the action of the lens at $z = z_L$ introduces a path difference of $\Delta l = [f - (f^2 - x_L^2)^{1/2}] \simeq (x_L^2/2f)$ at a transverse distance x_L . This corresponds to a phase difference $k\Delta l = (kx_L^2/2f) = \theta$ introduced by the lens.

Let us next consider the effect of this lens on a point source of radiation along the z axis at $z = z_O$. [That is, the initial amplitude is taken to be to be $a(z_O, \mathbf{x}_O) \propto \delta_D(\mathbf{x}_O)$.] This can be obtained by first propagating the field from z_O to z_L , modifying the phase due to the lens at $z = z_L$ and propagating it further to some point z with the transverse coordinate set to zero. The net result is given by

Does it really work as a lens?

$$a(z, 0) = -\frac{k^2}{4\pi^2 uv} \int d^2\mathbf{x}_L \exp\left(-\frac{ik}{2f}x_L^2\right) \cdot \exp\left[\frac{ikx_L^2}{2u} + \frac{ikx_L^2}{2v}\right], \quad (16.17)$$

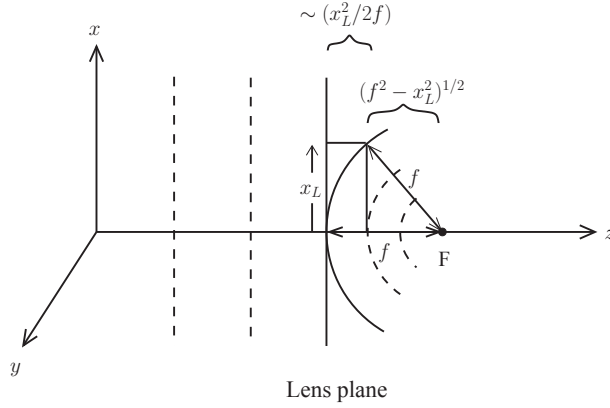


Fig. 16.1: The focusing action of a convex lens in terms of the phase change of wave fronts.

where $u = z_L - z_O$ and $v = z - z_L$. In the limit of zero wavelength (called *ray optics*), the maximum contribution to this integral can again be obtained by setting the variation of the phase to zero. This gives

It does!

$$-\frac{k}{f}\mathbf{x}_L + \frac{k}{u}\mathbf{x}_L + \frac{k}{v}\mathbf{x}_L = 0, \quad (16.18)$$

or

$$\frac{1}{u} + \frac{1}{v} = \frac{1}{f}, \quad (16.19)$$

which is a familiar formula in the theory of lenses.

Actually we can do better

The above result was obtained in the limit of ray optics. To study the wave propagation through the lens, we note that the action of a lens on the phase of an initial intensity distribution is governed by the integral

$$a(z, \mathbf{x}) \propto \int d^2 \mathbf{x}_L a(z_L, \mathbf{x}_L) \exp\left(-\frac{ik}{2f}x_L^2\right) \exp\frac{ik}{2(z-z_L)}(\mathbf{x}-\mathbf{x}_L)^2. \quad (16.20)$$

Here, $a(z_L, \mathbf{x}_L)$ is the incident amplitude on the lens; the first exponential gives the distortion in phase produced by the lens and the second exponential gives the propagation amplitude z_L to z . At the focal plane, which is a plane located at a distance f from the lens, at $z = z_L + f$, the second exponential characterizing the propagation becomes:

$$\exp\frac{ik(\mathbf{x}-\mathbf{x}_L)^2}{2(z-z_L)} = \exp\frac{ik}{2f}(x^2 + x_L^2 - 2\mathbf{x} \cdot \mathbf{x}_L). \quad (16.21)$$

The quadratic term $(ikx_L^2/2f)$ in the propagation amplitude is now precisely canceled by the phase distortion introduced by the lens, so that the resultant amplitude can be written as

$$a(z_L + f, \mathbf{x}) \propto \exp\left(\frac{ik}{2f}x^2\right) \int d^2\mathbf{x}_L a(z_L, \mathbf{x}_L) \exp\left(\frac{ik}{f}\mathbf{x} \cdot \mathbf{x}_L\right). \quad (16.22)$$

The intensity at the focal plane is given by the $|a(z_L + f, \mathbf{x})|^2$ in which the phase factor $\exp[ikx^2/2f]$ does not contribute. This is clearly determined by the Fourier transform of the incident amplitude. Thus we find that our humble lens acts as an analogue machine which performs the Fourier transform of a function.

Lens calculates a Fourier transform!

Box 16.1: Diffraction from Faraday's law

Why does light, treated as electromagnetic wave exhibit diffraction when it passes through a small aperture in a screen? In the standard approach, one first obtains the electromagnetic wave equation by combining the individual Maxwell equations suitably and then derives diffraction as a standard result in wave propagation. At this stage, the diffraction of light is no different from the diffraction of sound. But unlike sound, we know that the electromagnetic field has to satisfy *each* of the Maxwell equations separately. Using this fact, one can provide an intuitive understanding of diffraction at an aperture.

Why diffraction?

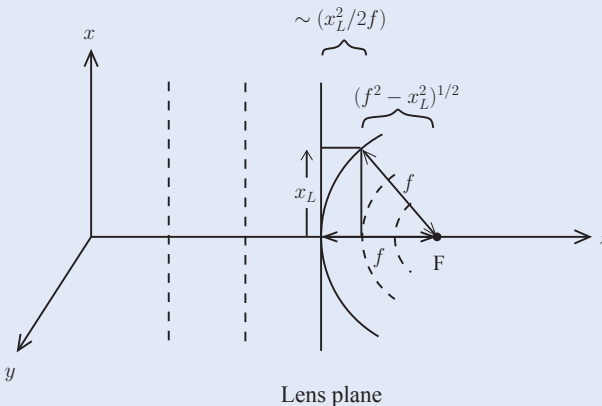


Fig. 16.2: A simple way to understand diffraction using Faraday's law. An electromagnetic wave propagates along z -axis, passing through a square aperture in a screen with the electric field along the x -axis and magnetic field along the y -axis originally. After passing through the aperture, the line integral of the electric field along the curve shown will be non-zero which requires a time varying z -component for the magnetic field. This requires the propagation direction to change slightly, which leads to the diffraction spread.

The crucial input

Consider a linearly polarized electromagnetic wave, with the electric field along the x -axis and propagating along the z -axis. Suppose this wave passes through a square aperture of size ℓ located in the $z = 0$ plane. We study the line integral of the electric field along the contour indicated in Fig. 16.2. This contour is parallel to the screen and is very close to it on the other side of the source. This line integral is essentially given by ℓE_x on the side away from the source. Numerically, this is also equal to ℓB_y . But Faraday told us that this must be equal to the rate of change of magnetic flux through the loop. In other words, we must have a z -component of the magnetic field generated on the far side even though none was present originally! Taking $\partial B_z / \partial t = -i\omega B_z$, we get the result

$$\ell B_y = \ell E_x = \oint \mathbf{E} \cdot d\mathbf{s} = -\frac{1}{c} \frac{\partial}{\partial t} \int \mathbf{B} \cdot d\mathbf{a} = -\frac{1}{c} (-i\omega B_z) \frac{\ell^2}{2}. \quad (16.23)$$

The first equality comes from $E_x = B_y$ for an electromagnetic wave, the second from the estimate of the line integral, the third from Faraday's law and the fourth from an estimate of the surface integral. We therefore get the longitudinal component of the magnetic field generated after the screen to be given by

$$\frac{B_z}{B_y} \cong -\frac{i\lambda}{\pi\ell}, \quad (16.24)$$

where λ is the wavelength of the radiation. This clearly gives the estimate of the standard diffraction angle to be about λ/ℓ . (The i -factor in the above relation also contains important information about the phase but we will not go into it here.)

If Quantum Mechanics is the Paraxial Optics, then ...

17

In quantum mechanics, the wavefunction of the particle, $\psi(t, \mathbf{x})$ contains complete information about the state of the system and satisfies the Schrödinger equation. Given the wavefunction $\psi(0, \mathbf{x})$ at $t = 0$, we can integrate this equation and obtain the wavefunction at any later time. So, all the dynamics is contained in the probability amplitude $\langle x_2 | x_1 \rangle$ for the particle to propagate from one event x_1 to another event x_2 . (For example, the beaten-to-death electron two slit experiment involves an electron gun to create electrons and a detector on the screen to detect them). Classically, the particle will move from one event x_1 to another event x_2 along a single, deterministic, trajectory. But we know that, in quantum mechanics, there is no notion of trajectories at all. Is there some nice way of expressing this *quantum* amplitude $\langle x_2 | x_1 \rangle$ in terms of what we know in *classical* physics?

Initial value problem in quantum mechanics

A hint that it may be possible arises from our results in Chapter 2 where we saw that the *classical action* plays a crucial role even in quantum theory and — in fact — it is quantum mechanics which validates the principle of least action in classical theory [67]. We could define an action for all possible trajectories by Eq. (2.30) and recover the classical trajectory through the condition for stationary phase. Since the classical action A and the quantum amplitude Ψ are related by $\Psi \propto \exp(iA/\hbar)$, it seems natural to postulate that the amplitude for a particle to follow a particular trajectory $\mathbf{x}(t)$ is proportional to $\exp(iA[\mathbf{x}(t)]/\hbar)$ where $A[\mathbf{x}(t)]$ is the action for that trajectory. This postulate assures us at least one thing: In the classical limit of $\hbar \rightarrow 0$, the condition for constructive interference will pick out the classical path! Since all paths are possible in the fully quantum mechanical situation, the net amplitude $\langle x_2 | x_1 \rangle$ for the particle to go from one event x_1 to another event x_2 must be the sum over $\exp(iA[\mathbf{x}(t)]/\hbar)$ for all paths connecting the two events. So, it seems natural to expect:

Action: the common factor

$$\langle x_2 | x_1 \rangle = \sum_{\mathbf{x}(t)} \exp \left[\frac{iA[\mathbf{x}(t)]}{\hbar} \right] = \sum_{\mathbf{x}(t)} \exp \frac{i}{\hbar} \int_{t_1}^{t_2} \frac{1}{2} m |\dot{\mathbf{x}}|^2 dt . \quad (17.1)$$

You don't really sum over all paths!

The paths summed over are restricted to those that satisfy the following condition: *Any given path $\mathbf{x}(t)$ cuts the spatial hypersurface $t = y^0$ at any intermediate time, $t_2 > y^0 > t_1$, at only one point.* In other words, while doing the sum over paths, we are restricting ourselves to paths of the kind shown in Fig. 17.1 that always go ‘forward in time’ and do not include, for example, paths like the one shown in Fig. 17.2 (which go both forward and backward in time).

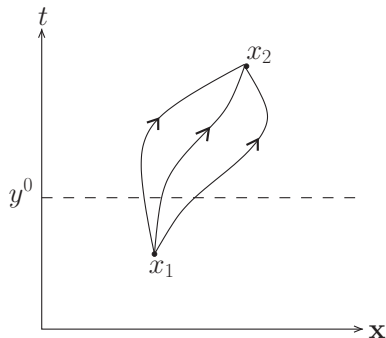


Fig. 17.1: Examples of paths included in the sum over paths in Eq. (17.1)

Why we don't want certain type of paths

The path in Fig. 17.2 cuts the constant time surface $t = y^0$ at three events, suggesting that at $t = y^0$ there were three particles simultaneously present even though we started out with one particle. It is this feature which we avoid (and stick to single particle propagation) by imposing this condition on the class of paths that is included in the sum. By the same token, we will assume that the amplitude $\langle x_2 | x_1 \rangle$ vanishes for $x_2^0 < x_1^0$; that is, the propagation is forward in time.

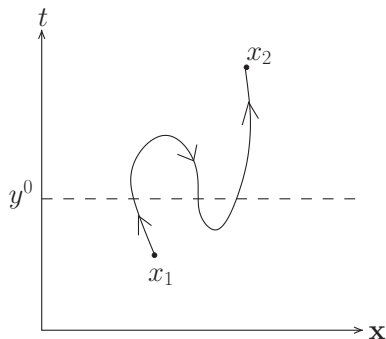


Fig. 17.2: A path that goes forward and backward which is *not* included in the sum over paths in Eq. (17.1).

This choice of paths, in turn, implies the following ‘transitivity constraint’ for the amplitude:

$$\langle x_2 | x_1 \rangle = \int d^D \mathbf{y} \langle x_2 | \mathbf{y} \rangle \langle \mathbf{y} | x_1 \rangle . \quad (17.2)$$

The integration at an intermediate event $y \equiv y^i = (y^0, \mathbf{y})$ (with $t_2 > y^0 > t_1$) is limited to integration over the *spatial* coordinates because each of the paths summed over cuts the intermediate spatial surface at only one point. Therefore, every path which connects the events x_2 and x_1 can be uniquely specified by the spatial location \mathbf{y} at which it crosses the surface $t = y^0$. So the sum over all paths can be divided into the sum over all paths from x_1^i to some location \mathbf{y} at $t = y^0$, followed by the sum over all paths from y^i to x_2^i with an integration over all the locations \mathbf{y} at the intermediate time $t = y^0$. This leads to Eq. (17.2).

A nontrivial demand ...

The transitivity condition in Eq. (17.2) is *vital* for the standard probabilistic interpretation of the wavefunction in non-relativistic quantum mechanics. If $\psi(t_1, \mathbf{x}_1)$ is the wavefunction giving the amplitude to find a particle at \mathbf{x}_1 at time t_1 , then the wavefunction at a later time $t = y^0$ is given by the integral:

... but quite essential

$$\psi(y^0, \mathbf{y}) = \int d^D \mathbf{x}_1 \langle \mathbf{y} | x_1 \rangle \psi(t_1, \mathbf{x}_1) , \quad (17.3)$$

which interprets $\langle \mathbf{y} | x_1 \rangle$ as a propagator kernel allowing us to determine the solution to a differential equation (viz. the Schrödinger equation) at a later time $t = y^0$ from its solution at $t = t_1$. Writing the expression for $\psi(t_2, \mathbf{x}_2)$ in terms of $\psi(y^0, \mathbf{y})$ and $\langle x_2 | \mathbf{y} \rangle$ and using Eq. (17.3) to express $\psi(y^0, \mathbf{y})$ in terms of $\psi(t_1, \mathbf{x}_1)$, it is easy to see that Eq. (17.2) is *needed* for consistency. Equation (17.2) or Eq. (17.3) also implies the condition:

$$\langle t, \mathbf{x} | t, \mathbf{y} \rangle = \delta(\mathbf{x} - \mathbf{y}) , \quad (17.4)$$

where $|t, \mathbf{x}\rangle$ is a position eigenstate at time t .

Three crucial factors have gone into these seemingly innocuous results:

(i) The wavefunction at time t can be obtained from knowing only the wavefunction at an earlier time (without, e.g., knowing its time derivative). This means that the differential equation governing ψ must be first order in time. (ii) One can introduce eigenstates $|t, \mathbf{x}\rangle$ of the position operator $\hat{\mathbf{x}}(t)$ at time t by $\hat{\mathbf{x}}(t)|t, \mathbf{x}\rangle = \mathbf{x}|t, \mathbf{x}\rangle$ so that $\psi(t, \mathbf{x}) = \langle t, \mathbf{x} | \psi \rangle$ with Eq. (17.4) allowing the possibility of localizing a particle in space with arbitrary accuracy. (iii) One can interpret $\langle x_2 | x_1 \rangle$ in terms of the position eigenstates as $\langle t_2, \mathbf{x}_2 | t_1, \mathbf{x}_1 \rangle$. It turns out all these conditions run into trouble when we deal with a relativistic particle! This is why quantum field theory is very different from single particle quantum mechanics.

Three assumptions in quantum mechanics: All invalid in QFT!

*Free particle from
first principles*

Before proceeding further, let us consider the case of a free-particle in this formalism. Here, we can make a lot of progress by just using the result that the integral in Eq. (17.2) has to be independent of y^0 , since the left hand side is independent of y^0 . The transitivity condition, plus the fact that the free particle amplitude $\langle x_2|x_1 \rangle$ can only depend on $|\mathbf{x}_2 - \mathbf{x}_1|$ and $(t_2 - t_1)$ because of translational and rotational invariance, fixes the form of $\langle x_2|x_1 \rangle$ to a great extent. To see this, express $\langle x_2|x_1 \rangle$ in terms of its spatial Fourier transform in the form

$$\langle y|x \rangle = \int \frac{d^D \mathbf{p}}{(2\pi)^D} \theta(y^0 - x^0) F(|\mathbf{p}|; y^0 - x^0) e^{i\mathbf{p} \cdot (\mathbf{y} - \mathbf{x})}, \quad (17.5)$$

and substitute into Eq. (17.2). This will lead to the condition

$$F(|\mathbf{p}|; x_2^0 - y^0) F(|\mathbf{p}|; y^0 - x_1^0) = F(|\mathbf{p}|; x_2^0 - x_1^0) \quad (x_2^0 > y^0 > x_1^0), \quad (17.6)$$

which has a unique solution $F(|\mathbf{p}|; t) = \exp(\alpha(|\mathbf{p}|)t)$ where $\alpha(|\mathbf{p}|)$ is a function of $|\mathbf{p}|$. Further, we note that $F(|\mathbf{p}|; y^0 - x^0)$ propagates the momentum space wavefunction $\phi(x^0, \mathbf{p})$ — which is the spatial Fourier transform of $\psi(x^0, \mathbf{x})$ — from $t = x^0$ to $t = y^0$. Since ϕ is the Fourier transform of ψ , this “propagation” is just multiplication by F . The probability calculated from the momentum space wavefunction will be well behaved for $|t| \rightarrow \infty$ only if α is pure imaginary, thereby only contributing a phase. So $\alpha = -if(|\mathbf{p}|)$ where $f(|\mathbf{p}|)$ is an arbitrary function of $|\mathbf{p}|$. (You can also obtain the same result from the fact that $\exp(iA)$ goes to $\exp(-iA)$ under the time reversal $t_2 \iff t_1$; the path integral sum must be defined such that $F \rightarrow F^*$ under $t \rightarrow -t$ which requires α to be pure imaginary.) Thus, the spatial Fourier transform of $\langle x_2|x_1 \rangle$ must have the form

$$\int d^D \mathbf{x} \langle x_2|x_1 \rangle e^{-i\mathbf{p} \cdot \mathbf{x}} = \theta(t) e^{-if(|\mathbf{p}|)t}. \quad (17.7)$$

That is, it must be a pure phase. If we interpret this phase as due to the energy $\omega_{\mathbf{p}} = \mathbf{p}^2/2m$ and set $f(\mathbf{p}) = \omega_{\mathbf{p}}$, then an inverse Fourier transform of Eq. (17.7) will immediately determine $\langle x_2|x_1 \rangle$ leading to the result:

$$\begin{aligned} \langle x_2|x_1 \rangle &\equiv K(t, \mathbf{x}; 0, \mathbf{y}) = \int \frac{d^D \mathbf{p}}{(2\pi)^D} e^{i\mathbf{p} \cdot (\mathbf{x} - \mathbf{y})} e^{-ip^2 t/2m} \\ &= \left(\frac{m}{2\pi i \hbar t} \right)^{D/2} \exp \left[\frac{im}{\hbar} \frac{(\mathbf{x} - \mathbf{y})^2}{2t} \right], \end{aligned} \quad (17.8)$$

where D is the dimension of space (1, 2 or 3) in which the particle is moving and we have reintroduced the \hbar . The integral is just the D -dimensional Fourier transform of a Gaussian which separates out in each of the dimensions. We can verify directly that $K(t, \mathbf{x}; 0, \mathbf{y})$ satisfies Eq. (17.2) and Eq. (17.4).

Almost there

*Just one more
input ...*

*... and we have the
final result*

The sum over paths in Eq. (17.1) itself is trivial to evaluate for all classical actions, which are at most quadratic in $\mathbf{x}(t)$ and $\dot{\mathbf{x}}(t)$, even without us defining precisely what the sum means. (The more sophisticated definitions for the sum work — or rather designed to work — only because we know the answer for $\langle x_2|x_1 \rangle$ from other well-founded methods!) We first note that the sum over all $\mathbf{x}(t)$ is the same as the sum over all $\mathbf{q}(t) \equiv \mathbf{x}(t) - \mathbf{x}_c(t)$ where $\mathbf{x}_c(t)$ is the classical path for which the action is an extremum. Because of the extremum condition, $A[\mathbf{x}_c + \mathbf{q}] = A[\mathbf{x}_c] + A[\mathbf{q}]$. Substituting into Eq. (17.1) and noting that $\mathbf{q}(t)$ vanishes at the end points, we see that the sum over $\mathbf{q}(t)$ must be only a function of $(t_2 - t_1)$. (It can only depend on the time difference rather than on t_2 and t_1 individually whenever the action has no explicit time dependence; i.e., for any closed system). Thus we get

We can do better but not a lot better

$$\langle x_2|x_1 \rangle = e^{iA[\mathbf{x}_c]} \sum_{\mathbf{q}} e^{iA[\mathbf{q}(t)]} = N(t) \exp iA[\mathbf{x}_c] , \quad (17.9)$$

where $t \equiv t_2 - t_1$. Thus, the *quantum* probability amplitude is expressible in terms of the *classical* action for the *classical* trajectory, except for a normalization function $N(t)$, for all quadratic actions. This factor needs to be determined by some other clever trick in each case. For the free particle, we immediately get:

This is useful

$$\langle x_2|x_1 \rangle = N(t) \exp iA[\mathbf{x}_c] \equiv N(t) \exp \left(\frac{i}{2} \frac{m|\mathbf{x}|^2}{t} \right) , \quad (17.10)$$

where $t \equiv t_2 - t_1$, $\mathbf{x} \equiv \mathbf{x}_2 - \mathbf{x}_1$ and $\hbar = 1$. In this case, the form of $N(t)$ is strongly constrained by the transitivity condition, Eq. (17.2) — or, equivalently, by Eq. (17.7) — which requires the $N(t)$ to have the form $(m/2\pi it)^{D/2} e^{a t}$ where $a = i\phi$, say. Thus, except for an ignorable, constant, phase factor ϕ (which is equivalent to adding a constant to the Lagrangian), $N(t)$ is given by $(m/2\pi it)^{D/2}$ and we can write the full propagation amplitude for a non-relativistic particle as:

$$\langle x_2|x_1 \rangle = \theta(t) \left(\frac{m}{2\pi it} \right)^{D/2} \exp \left(\frac{i}{2} \frac{m|\mathbf{x}|^2}{t} \right) . \quad (17.11)$$

The $\theta(t)$ tells you that we are considering a particle which is created, say, at t_1 and detected at t_2 with $t_2 > t_1$. In non-relativistic mechanics, all inertial observers will give an invariant meaning to the statement $t_2 > t_1$. It is also easy to see that the $\langle x_2|x_1 \rangle$ in Eq. (17.11) satisfies the condition in Eq. (17.4).

I said that we can compute the path integral only for quadratic actions. This is by and large true but there is one peculiar (and important) case of a non-quadratic action for which the path integral can be evaluated exactly by a trick. Given the fact that it is not as well-known as it should be, let

*A non-quadratic
path integral*

me describe this. The trick here uses the fact that the action functional for a particle in classical mechanics can also be expressed in the Jacobi-Mapertuis form (discussed in Chapter 2) which has a square root in it. We saw that the trajectory of the particle can be obtained in classical theory from the action expressed in the form (see Eq. (2.39)):

$$A_J = \int_{x_1}^{x_2} m \left(\frac{dl}{d\lambda} \right) d\lambda = \int_{x_1}^{x_2} \sqrt{2m(E - V(x^\alpha))} d\lambda. \quad (17.12)$$

Since A_J describes a valid action principle for finding the path of a particle with energy E classically, one might wonder what happens if we try to quantize the system by performing a sum over amplitudes $\exp(iA_J)$. We would expect it to lead to the amplitude for the particle to propagate from x_1^α to x_2^α with energy E . This is indeed true, but since A_J is not quadratic in velocities even for a free particle, (note that dl involves a square root) it is not easy to evaluate the sum over $\exp(iA_J)$. But since we already have an alternative path integral procedure for the system, we can use it *to give meaning* to this sum, thereby evaluating the sum over paths for at least one non-quadratic action.

Our idea is to write the sum over all paths in the original action principle (with amplitude $\exp(iA)$) as a sum over paths with energy E followed by a sum over all E . Using the result in Eq. (2.10), we get

$$\sum_{0, x_1}^{t, x_2} \exp(iA) = \sum_E \sum_{x_1}^{x_2} e^{-iEt} \exp(iA_J[E, \mathbf{x}(\tau)]) \propto \int_0^\infty dE e^{-iEt} \sum_{x_1}^{x_2} \exp(iA_J). \quad (17.13)$$

In the last step, we have treated the sum over E as an integral over $E \geq 0$ (since, for any Hamiltonian which is bounded from below, we can always achieve this by adding a suitable constant to the Hamiltonian) but there could be an extra proportionality constant which we cannot rule out. This constant will depend on the measure used to define the sum over $\exp(iA_J)$ but can be fixed by using the known form of the left hand side, if required. Inverting the Fourier transform, we get:

Useful result

$$\begin{aligned} \mathcal{P}(E; \mathbf{x}_2, \mathbf{x}_1) &\equiv \sum_{x_1}^{x_2} \exp(iA_J) = C \int_0^\infty dt e^{iEt} \sum_{0, x_1}^{t, x_2} \exp(iA) \\ &= C \int_0^\infty dt e^{iEt} \langle x_2 | x_1 \rangle, \end{aligned} \quad (17.14)$$

where we have denoted the proportionality constant by C . This result shows that the sum over the Jacobi action involving a *square root of velocities* can be re-expressed in terms of the standard path integral; if the latter can be evaluated for a given system, then the sum over the Jacobi action can be defined by this procedure.

The result also has an obvious interpretation. The $\langle x_2|x_1 \rangle$ on the right hand side gives the amplitude for a particle to propagate from \mathbf{x}_1 to \mathbf{x}_2 in time t . Its Fourier transform with respect to t can be thought of as the amplitude for the particle to propagate from \mathbf{x}_1 to \mathbf{x}_2 with energy E , which is precisely what we expect to obtain from the sum over the Jacobi action. The idea actually works even for particles in a potential if we evaluate the path integral on the right hand side by some other means like, e.g., by solving the relevant Schrödinger equation.

With future applications in mind, we will display the explicit form of this result for the case of a free particle with $V = 0$. Denoting the length of the path connecting x_1^α and x_2^α by $\ell(\mathbf{x}_2, \mathbf{x}_1)$ we have:

$$\sum_{\mathbf{x}_1}^{\mathbf{x}_2} \exp i\sqrt{2mE} \ell(\mathbf{x}_2, \mathbf{x}_1) = C \int_0^\infty dt e^{iEt} \sum_{0, \mathbf{x}_1}^{\mathbf{x}_2} \exp \frac{im}{2} \int_0^t d\tau \left(g_{\alpha\beta} \dot{x}^\alpha \dot{x}^\beta \right). \quad (17.15)$$

This result shows that the sum over paths with a Jacobi action, which has a square root, can be re-expressed in terms of the standard path integral involving only quadratic terms in the velocities. We, of course, know the result of the path integral in the right hand side (for $g_{\alpha\beta} = \delta_{\alpha\beta}$ in Cartesian coordinates) and thus we can evaluate the sum on the left hand side.

Box 17.1: Propagation amplitude from Stationary states

For a particle in a general, non-quadratic potential $V(x)$, nobody knows how to sum over paths and get $\langle x_2|x_1 \rangle$. So the path integral is nice to look at but practically useless — you need to get back to the Schrödinger equation! But one can certainly express $\langle x_2|x_1 \rangle$ in terms of solutions to the Schrödinger equation, when the potential is time-independent, as follows:

When the potential is independent of time, energy eigenstates satisfy the eigenvalue equation $H\phi_n(\mathbf{x}) = E_n\phi_n(\mathbf{x})$. Using these eigenfunctions we can expand the initial wavefunction $\psi(0, \mathbf{x})$ in terms of the energy eigenfunctions as

$$\psi(0, \mathbf{x}) = \sum_n c_n \phi_n(\mathbf{x}); \quad c_n = \int d\mathbf{y} \psi(0, \mathbf{y}) \phi_n^*(\mathbf{y}), \quad (17.16)$$

where the expression for c_n follows from the orthonormality of the energy eigenfunctions and the spatial integrations are over the D -dimensional space. Since the energy eigenfunction evolves in time with a phase factor $\exp(-iE_n t/\hbar)$, it follows that the wavefunction at time t is given by:

$$\psi(t, \mathbf{x}) = \sum_n c_n \phi_n(\mathbf{x}) e^{-iE_n t/\hbar}, \quad (17.17)$$

*Come back,
Schrödinger; all
is forgiven!*

which, in principle, solves the problem. We now express the c_n s in Eq. (17.17) in terms of $\psi(0, \mathbf{x})$ using the second relation in Eq. (17.16). This gives:

$$\begin{aligned}\psi(t, \mathbf{x}) &= \int d\mathbf{y} \psi(0, \mathbf{y}) \sum_n \phi_n(\mathbf{x}) \phi_n^*(\mathbf{y}) e^{-iE_n t/\hbar} \\ &\equiv \int d\mathbf{y} K(t, \mathbf{x}; 0, \mathbf{y}) \psi(0, \mathbf{y}) ,\end{aligned}\quad (17.18)$$

which allows us to read off the propagator as:

$$\langle x_2 | x_1 \rangle \equiv K(t, \mathbf{x}; 0, \mathbf{y}) = \sum_n \phi_n^*(\mathbf{y}) \phi_n(\mathbf{x}) e^{-iE_n t/\hbar} . \quad (17.19)$$

Equation (17.18) nicely separates the dynamics — encoded in $K(t, \mathbf{x}; 0, \mathbf{y})$ — from the initial condition encoded in $\psi(0, \mathbf{y})$. *Curiously enough, such a separation has no direct analog in the case of classical mechanics.*

Using the definition in Eq. (17.19) and the orthonormality of eigenfunctions, you can prove that $\langle x_2 | x_1 \rangle$ does satisfy the two constraints in Eq. (17.2) and Eq. (17.4). Since the ϕ_n s are energy eigenfunctions, it is also straightforward to verify that the propagator satisfies the Schrödinger equation

$$\left(i\hbar \frac{\partial}{\partial t} - H \right) K(t, \mathbf{x}; 0, \mathbf{y}) = 0 , \quad (17.20)$$

with the special initial condition

$$\lim_{t \rightarrow 0} K(t, \mathbf{x}; 0, \mathbf{y}) = \delta_D(\mathbf{x} - \mathbf{y}) . \quad (17.21)$$

This condition can also be obtained easily from Eq. (17.18) by taking the limit of $t \rightarrow 0$.

We said earlier that the exact evaluation of the sum over paths is possible only when the action is quadratic. But, there are situations in which one can *approximate* the sum by the result in Eq. (17.9) which only requires the *classical* solution to the problem. Then, by using Eq. (17.19), we can get some information about energy spectrum — which is (approximate) quantum mechanics at the classical price! We will say more about this in Chapter 18.

The most remarkable feature about the propagator in Eq. (17.11) is that *you have already seen this expression* in Chapter 16 in connection with the propagation of electromagnetic waves along the z -direction! There we had the expression (see Eq. (16.8)) for a propagator which is reproduced here

*Propagator
from energy
eigenfunctions*

A strange fact

*Here comes the
real surprise*

for your convenience:

$$G(z-z'; \mathbf{x}_\perp - \mathbf{x}'_\perp) = \left(\frac{\omega}{2\pi i c} \right) \frac{1}{|z-z'|} \exp \left[\frac{i\omega}{2c} \frac{(\mathbf{x}_\perp - \mathbf{x}'_\perp)^2}{(z-z')} \right]. \quad (17.22)$$

Comparing Eq. (17.22) with Eq. (17.8), we see the following correspondence. The $(z-z')/c$, which is the time of light travel along the z -axis — along which the wave is propagating — is analogous to the time t in quantum mechanics. The two transverse spatial directions in the case of electromagnetic wave propagation are analogous to the spatial coordinates in quantum mechanics in 2-dimensions, so that we can set $D = 2$ in Eq. (17.8). The frequency should get mapped to the relation $\hbar\omega = mc^2$ which is essentially the frequency associated with the Compton wavelength of the particle. This will make the propagators identical! Obviously, this deserves further probing especially since the correspondence brings in a c factor when we thought we were doing non-relativistic quantum mechanics.

Where did c spring from in quantum mechanics?!

In the case of the propagation of the electromagnetic wave amplitude, we were propagating it along the positive z -direction with \mathbf{x} and \mathbf{y} acting as two transverse directions. In the case of quantum mechanics, we are propagating the amplitude for a particle along the positive t -direction with all the spatial coordinates acting as “transverse directions”. In the language of paraxial optics, the special axis is along the *time direction* in quantum mechanics.

Quantum mechanics is paraxial optics in time direction! You go forward but not backward in time!

But we know that paraxial optics is just an approximation to a more exact propagation in terms of the wave equation. In the wave equation for the electromagnetic wave, the three coordinates (x, y, z) appear quite symmetrically and to obtain the paraxial limit, we choose one axis (the z -axis) as special and propagate the amplitude along the positive direction. This is why the propagator in Eq. (17.22) has the x, y coordinates appearing differently compared to the z -axis. Doing a bit of reverse engineering we can ask the question: *If the quantum mechanical propagator is some kind of paraxial optics limit of a more exact theory, what is the exact theory?*

If so, what is the real thing?

An obvious way to explore the situation is to restore the symmetry between z and x, y in optics and — similarly — restore the symmetry between t and \mathbf{x} in quantum mechanics. We can do this if we recall the interpretation of the phase as due to the path difference in the case of an electromagnetic wave. The relevant equation (see Eq. (16.10)) is again reproduced below:

$$\begin{aligned} k\Delta s &= \frac{\omega}{c} \left[\sqrt{(\mathbf{x}_\perp - \mathbf{x}'_\perp)^2 + (z-z')^2} - (z-z') \right] \\ &\cong \frac{\omega}{c} \left[\frac{1}{2} \frac{(\mathbf{x}_\perp - \mathbf{x}'_\perp)^2}{(z-z')} \right]. \end{aligned} \quad (17.23)$$

Ha! from quantum mechanics to better things!

We use the fact that a path difference Δs between two points in space will introduce a phase difference of $k\Delta s$ in a propagating wave. The paraxial optics results when the transverse displacements are small compared to the longitudinal distance. Taking a cue from this, let us construct the quantity

$$\frac{\ell(t, \mathbf{x}; 0, \mathbf{y}) - ct}{\lambda} \equiv \frac{mc}{\hbar} \left\{ [c^2 t^2 - (\mathbf{x} - \mathbf{y})^2]^{1/2} - ct \right\}, \quad (17.24)$$

where $\ell(t, \mathbf{x}; 0, \mathbf{y})$ is the special relativistic spacetime interval between the two events. We are subtracting from it the “paraxial distance” ct along the time direction and dividing by $\lambda \equiv (\hbar/mc)$ which is the Compton wavelength of the particle. This is exactly the construction suggested by the correspondence between Eq. (17.22) and Eq. (17.8), discussed previously, except for using the special relativistic line interval, with a minus sign between space and time. The paraxial limit now arises as the non-relativistic limit of this expression in Eq. (17.24) when $c \rightarrow \infty$; this is given by:

$$\frac{\ell - ct}{\lambda} \cong -\frac{m(\mathbf{x} - \mathbf{y})^2}{2\hbar t}, \quad (17.25)$$

which is precisely the phase of the propagator in Eq. (17.8) except for a sign. So, the propagator can be thought of as the non-relativistic limit of the function:

$$K(t, \mathbf{x}; 0, \mathbf{y}) = N(t) e^{i(mc^2/\hbar)t} \exp \left(-i \left[\frac{\ell(t, \mathbf{x}; 0, \mathbf{y})}{\lambda} \right] \right). \quad (17.26)$$

So, the phase of the propagator is just the proper distance between the two events, in units of the Compton wavelength, just as the phase in the case of the electromagnetic wave propagator is the path length in units of the wavelength. (The extra factor $(mc^2/\hbar)t$ does not contribute to the propagation integral in Eq. (17.18) and goes for a ride in this context; however, it has some curious implications which we discussed in Chapter 15.). We can think of the path difference between a straight path along the time direction (with $\mathbf{x} = \mathbf{y}$) and another specified path as contributing a phase ℓ/λ to the propagator. This geometric interpretation is lost for the phase in the paraxial limit (in the case of electromagnetic theory) and in the non-relativistic limit (in the case of a particle).

Approximate theories take away all the fun

This extension suggests that the phase in the relativistic case can be related to the corresponding action. The action for a free particle in special relativity is given by

$$A_R(t, \mathbf{x}; 0, \mathbf{y}) = -mc^2 \int_0^t dt \left(1 - \frac{v^2}{c^2} \right)^{1/2}. \quad (17.27)$$

Once again, evaluating this for a relativistic classical trajectory, we get:

$$A_R^c(t, \mathbf{x}; 0, \mathbf{y}) = -mc^2 t \left[1 - \frac{(\mathbf{x} - \mathbf{y})^2}{c^2 t^2} \right]^{1/2} = -mc [c^2 t^2 - (\mathbf{x} - \mathbf{y})^2]^{1/2}, \quad (17.28)$$

which is essentially the interval between the two events in the spacetime. This suggests expressing the propagator for the relativistic free particle in the form:

$$K(t, \mathbf{x}; 0, \mathbf{y}) = N(t) \exp \left(\frac{iA_R^c}{\hbar} + \frac{imc^2 t}{\hbar} \right). \quad (17.29)$$

Path length is natural in special relativity

This result is true but only in an approximate sense, to the leading order; the actual propagator for a particle in relativistic quantum theory turns out to be more complicated. This is because the action in Eq. (17.27) for the relativistic particle is not quadratic and our previous result in Eq. (17.9) does not hold. But, to the leading order, all of it hangs together very nicely. The phase of the propagator is indeed the value of the classical action divided by \hbar and it is also given by the ratio of the spacetime interval between the events and the Compton wavelength. It is the second interpretation which makes the contact with optics so clear and is lacking when we do non-relativistic quantum mechanics.

Minor caveat

There is actually a valid mathematical reason for this to happen, which can be described qualitatively as follows: The Schrödinger equation describing the non-relativistic particle involves the first derivative with respect to time but the second derivative with respect to spatial coordinates. This works in non-relativistic mechanics in which time is special and absolute. In contrast, in relativistic theories, we treat time and space at a more symmetric footing and use a wave equation in which the second derivative with respect to time also appears. The solutions to such an equation will allow propagation of amplitudes both forward and backward in the a time coordinate just as it allows propagation both forwards and backwards in spatial coordinates. *When one takes the non-relativistic limit of the field theory, we select out the modes which only propagate forward in time.*

Make sure you understand this

This is exactly in analogy with paraxial optics we studied in Chapter 16. The basic equation for an electromagnetic wave will allow propagation in both the positive z -direction as well as the negative z -direction. But, when we consider a specific context of paraxial optics (for example, a beam of light hitting a couple of slits in a screen and forming an interference pattern, or light propagating through a lens and getting focused), we select out the modes which are propagating in the positive z -direction. It is therefore no wonder that the propagator in non-relativistic quantum mechanics is mathematically identical to that in paraxial optics!

*Let us do it in full
glory, for once*

Finally, just to whet your curiosity, let me describe the structure of the exact relativistic propagator for a free particle. We will use units with $c = 1$ in what follows.

The standard action for a relativistic particle is given by

$$\begin{aligned} A &= -m \int_{t_1}^{t_2} dt \sqrt{1 - \mathbf{v}^2} = -m \int_{x_1}^{x_2} \sqrt{-\eta_{ab} dx^a dx^b} \\ &= -m \int_{\lambda_1}^{\lambda_2} d\lambda \sqrt{-\eta_{ab} \dot{x}^a \dot{x}^b}, \end{aligned} \quad (17.30)$$

where $x^a(\lambda)$ gives a parameterized curve connecting the events x_1 and x_2 in the spacetime with the parameter λ . In the second and third forms of the expression, the integral is evaluated for any curve connecting the two events with limits of integration depending on the nature of the parametrization. (For example, we have chosen $x(\lambda = \lambda_1) = x_1$, $x(\lambda = \lambda_2) = x_2$, but the numerical value of the integral is independent of the parametrization and depends only on the curve. If we choose to use $\lambda = t$ as the parameter, then we reproduce the first expression from the last.) It is obvious that this action has the same structure as the Jacobi action for a free particle discussed in the last section.

To obtain the propagation amplitude $\langle x_2 | x_1 \rangle$ we need to do the path integral using the above action,

$$\begin{aligned} \langle x_2 | x_1 \rangle &= \sum_{0, \mathbf{x}_1}^{t, \mathbf{x}_2} \exp \left[-im \int_{t_1}^{t_2} dt \sqrt{1 - \mathbf{v}^2} \right] \\ &= \sum_{0, \mathbf{x}_1}^{t, \mathbf{x}_2} \exp \left[-im \int_0^\tau d\lambda \sqrt{i^2 - \dot{\mathbf{x}}^2} \right], \end{aligned} \quad (17.31)$$

which can be accomplished using the results obtained earlier [68]. We first take the complex conjugate of Eq. (17.15) (in order to get the overall minus sign in the action in Eq. (17.30)) and generalize the result from space to spacetime, leading to:

$$\sum_{\mathbf{x}_1}^{x_2} \exp -i\sqrt{2mE} l(\mathbf{x}_2, \mathbf{x}_1) = C \int_0^\infty d\tau e^{-iE\tau} \sum_{0, \mathbf{x}_1}^{t, \mathbf{x}_2} \exp -\frac{im}{2} \int_0^\tau d\lambda \left(-g_{ab} \dot{x}^a \dot{x}^b \right). \quad (17.32)$$

In order to get $-iml(\mathbf{x}_2, \mathbf{x}_1)$ on the left hand side we take $E = m/2$; i.e., we use the above formula with the replacements

$$\begin{aligned} E &= \frac{m}{2}; \quad g_{ab} = \text{dia}(-1, +1, +1, +1); \\ l &= \int_0^\tau \sqrt{-g_{ab} \dot{x}^a \dot{x}^b} d\lambda = \int_{t_1}^{t_2} dt \sqrt{1 - \mathbf{v}^2}. \end{aligned} \quad (17.33)$$

The path integral over the quadratic action can be immediately borrowed from Eq. (17.11) with $D = 4$, taking due care of the fact that in the quadratic action the i^2 enters with negative sign while $\dot{\mathbf{x}}^2$ enters with the usual positive sign. This gives an extra factor i to N and the answer is:

$$\langle x_2, \tau | x_1, 0 \rangle = \theta(\tau) i \left(\frac{m}{2\pi i \tau} \right)^2 \exp \left[\frac{i}{2} \frac{m x^2}{\tau} \right]. \quad (17.34)$$

The $\theta(\tau)$ is introduced for the same reason as $\theta(t)$ in Eq. (17.11) but will turn out to be irrelevant since we will integrate over it. Therefore the path integral we need to compute is given by

$$\begin{aligned} \langle x_2 | x_1 \rangle &= \sum_{0, \mathbf{x}_1}^{t, \mathbf{x}_2} \exp -im \int_{t_1}^{t_2} dt \sqrt{1 - \mathbf{v}^2} \\ &= C \int_0^\infty d\tau e^{-\frac{im\tau}{2}} i \left(\frac{m}{2\pi i \tau} \right)^2 \exp \left(\frac{im}{2\tau} x^2 \right) \\ &= (2Cm)(-i) \frac{m}{16\pi^2} \int_0^\infty \frac{ds}{s^2} e^{-ims} \exp \left(\frac{im}{4s} x^2 \right); \quad \tau = 2s, \end{aligned} \quad (17.35)$$

where C is a proportionality constant. We have thus given meaning to the sum over paths for the relativistic particle thereby obtaining $\langle x_2 | x_1 \rangle$. The integral expression also gives a nice interpretation for $\langle x_2 | x_1 \rangle$ which we will first describe before discussing this result.

The trajectory of a classical relativistic particle in spacetime is given by the four functions $x^i(\tau)$ where τ could be taken as the proper time shown by a clock which moves with the particle. (To be precise, this is one physically meaningful choice for timelike curves; for spacelike and null curves, the corresponding choices are proper length and what is known as the affine parameter.). Such a description treats space and time on an equal footing with $\mathbf{x}(\tau)$ and $t(\tau)$ being *dependent* variables and τ being the *independent* variable having an observer independent, absolute, meaning. This is a natural generalization of $\mathbf{x}(t)$ in non-relativistic mechanics with (x, y, z) being dependent variables and t being the independent variable having an observer independent, absolute, status. Let us now consider an action for the relativistic particle in the form

*A quadratic action
for relativistic
particle*

$$A[x(\tau)] = \frac{1}{4} m \int_0^s d\tau \dot{x}_a \dot{x}^a, \quad (17.36)$$

where $\dot{x}^a \equiv (dx^a/d\tau)$, etc. This action, of course, gives the correct equations of motion $d^2 x^a / d\tau^2 = 0$, but the overall constant in front of the integral — which is arbitrary as far as the classical equations of motion go — is chosen with some foresight.

Evaluating a path integral with this action will now lead to an amplitude of the form $\langle x_2, s | x_1, 0 \rangle$ which describes a particle propagating from an

We don't care about proper time lapse

event x_1 to an event x_2 when the proper time lapse is given by s . But we are interested in the amplitude $\langle x_2|x_1 \rangle$ and don't care what is the amount of proper time that has elapsed. Therefore we need to also sum over (i.e., integrate) all the proper time lapses with some suitable measure. Since the rest energy of the particle $mc^2 = m$ is conjugate to the proper time (which measures the lapse of time in the instantaneous co-moving Lorentz frame of the particle) it seems reasonable to choose this measure to be proportional to a phase factor e^{-ims} . Thus we have the relation

$$\langle x_2|x_1 \rangle = C_m \int_{-\infty}^{\infty} ds e^{-ims} \langle x_2, s|x_1, 0 \rangle = C_m \int_{-\infty}^{\infty} ds e^{-ims} \sum_{x(\tau)} e^{iA[x(\tau)]}, \quad (17.37)$$

This is just a convention

where C_m is a normalization constant possibly dependent on m , which we will fix later. (The amplitude $\langle x_2|x_1 \rangle$ in Eq. (17.11) has the dimensions of $(\text{length})^{-D}$, as it should. So, the $\langle x_2|x_1 \rangle$ in Eq. (17.37) will have the dimension $(\text{length})^{-3}$ after integrating over s , if C_m is dimensionless. People like it to have the dimensions of $(\text{length})^{-2}$ which is achieved by taking $C_m \propto (1/m)$.) We have kept the integration limits on s to be the entire real line but it will get limited to $(0, \infty)$ because of the $\theta(s)$ in the path integral. In the second equality, we have used the standard path integral prescription.

Always go forward, in proper time ...

Exactly as before, the sum over paths is now to be evaluated limiting ourselves to paths $x^i(\tau)$ which only go forward in the *proper time* τ just as the paths in Eq. (17.10) were limited to those which go forward in the Newtonian absolute time t . However, we now *have* to allow paths like the one shown in Fig. 17.2 which go back and forth in time t just as we allowed in Eq. (17.10) the paths which went back and forth in the y coordinate, say. The time coordinate $t(\tau)$ of a path now has the same status as the spatial coordinate, say $y(\tau)$, in the non-relativistic description. The special role played by the absolute Newtonian time t is taken over by the proper time τ in this description. This has important implications which we will come back to later on.

... but backwards as well in coordinate time

Since the action is now quadratic, the calculation is straightforward and we get:

$$\begin{aligned} \langle x_2|x_1 \rangle &= -(2Cm)i \left(\frac{m}{16\pi^2} \right) \int_0^\infty \frac{ds}{s^2} \exp \left(-ims + \frac{i}{4} \frac{mx^2}{s} \right) \\ &= -\frac{i}{16\pi^2} \int_0^\infty \frac{d\mu}{\mu^2} \exp \left(-i(m^2 - i\varepsilon)\mu + \frac{i}{4} \frac{x^2}{\mu} \right), \quad (17.38) \end{aligned}$$

where we have made three modifications to arrive at the second line. First, we have rescaled the variable s to μ by $s \equiv m\mu$. Second, we have made the choice $C = 1/2m$ which, as we shall, see matches with conventional results later on and — more importantly — allows us to take the $m \rightarrow 0$ limit, if we want to study zero mass particles. Finally, we have replaced

m^2 by $(m^2 - i\varepsilon)$, where ε is an infinitesimal positive constant, in order to make the integral convergent in the upper limit. This is, of course, the same result obtained earlier. The integral can be expressed in terms of the MacDonald function:

$$\langle x_2 | x_1 \rangle = \frac{m}{4\pi^2 i \sqrt{-x^2}} K_1(im\sqrt{-x^2}), \quad (17.39)$$

where, of course, $x^2 = -t^2 + |\mathbf{x}|^2$ and hence the square-root of $-x^2$ is imaginary for space-like intervals. However, the Fourier transform of $\langle x_2 | x_1 \rangle$ is a more tractable object:

$$\begin{aligned} \int \langle x_2 | x_1 \rangle e^{ip \cdot x} d^4x &= -\frac{i}{16\pi^2} \int_0^\infty \frac{d\mu}{\mu^2} e^{-i(m^2 - i\varepsilon)\mu} \int d^4x \exp\left(\frac{i}{4} \frac{x^2}{\mu} + ip \cdot x\right) \\ &= -\frac{i}{(p^2 + m^2 - i\varepsilon)}, \end{aligned} \quad (17.40)$$

and is used extensively in field theory.

Let us now consider the nature of paths we summed over to get this result like the one in Fig. 17.2. This has the crucial implication that, at some intermediate coordinate time y^0 , we have to consider a situation with 3 particles at 3 different locations in space! Besides, the particle is traveling backwards in coordinate time for part of the path! This is disturbing to someone who is accustomed to sensible physical evolution which proceeds monotonously forward in time from $t_1 < t_2 < t_3 \dots$ and hence it would be nice if we can reinterpret $\langle x_2 | x_1 \rangle$ in such a nice, causal manner. Let us see what is needed for this.

If we say that a single particle has three degrees of freedom (in $D=3$), then we start and end (at x_1 and x_2) with three degrees of freedom in Fig. 17.2. But if a path cuts a spatial slice at an intermediate time y^0 at k points (the figure is drawn for $k=3$), then we need to be able to describe $3k$ degrees of freedom at this intermediate time. Since k can be arbitrarily large, we conclude that if we want a description in terms of causal evolution going from t_1 to t_2 , then we need to use a mathematical description involving an *infinite number* of degrees of freedom. In the properly constructed field theory, the parts of the particle trajectory which are going back in coordinate time are interpreted as the trajectory of an antiparticle going forward in time.

*Fortunately,
nobody ever uses
this because ...*

*... it is much
simpler in
Fourier space*

*Oops! You start
with one particle
and get three ...*

*... actually an in-
finite number of
them!*

In Chapter 17, we discussed how one can study the time evolution of a quantum wavefunction using a path integral propagator expressed as a sum over paths. We also showed that, when the Hamiltonian H is time independent, the kernel can be expressed in terms of the energy eigenfunctions through the formula:

See Eq. (17.19)

$$K(T, q_2; 0, q_1) = \sum_n \psi_n(q_2) \psi_n^*(q_1) \exp(-iE_n T) . \quad (18.1)$$

So, if the energy eigenfunctions and eigenvalues are given, one can determine the kernel. (We will use the terms propagator and kernel interchangeably.)

There are, however, occasions in which one may be able to determine the kernel directly by evaluating or approximating the path integral. The question arises as to whether one can determine the energy eigenfunctions and eigenvalues by “inverting” the above relation. In particular, one is often interested in the ground state eigenfunction and the ground state energy of the system. Can one find this if the kernel is known?

Can the sum over paths tell us something really useful?

It can be done using an interesting trick [69] which very often turns out to be more than just a trick, having a rather perplexing domain of validity. To achieve this, let us do the unimaginable and assume that time is actually a complex quantity. We then analytically continue from the real values of time t to purely imaginary values $\tau = it$. In special relativity such an analytic continuation will change the line interval from Lorentzian to Euclidean form through

Make time imaginary!

$$ds^2 = -dt^2 + d\mathbf{x}^2 \rightarrow d\tau^2 + d\mathbf{x}^2 . \quad (18.2)$$

Because of this, one often calls quantities evaluated with analytic continuation to imaginary values of time as “Euclidean” quantities and denotes them with a subscript E (which should not be confused with energy!). If we now do the analytic continuation of the kernel in Eq. (18.1), we get the

result

$$K_E(T_E, q_2; 0, q_1) = \sum_n \psi_n(q_2) \psi_n^*(q_1) \exp(-E_n T_E) . \quad (18.3)$$

Let us consider the form of this expression in the limit of $T_E \rightarrow \infty$. If the energy eigenvalues are ordered as $E_0 < E_1 < \dots$, then, in this limit, only the term with the ground state energy will make the dominant contribution, and remembering that ground state wavefunction is real for the systems we are interested in, we get,

$$K_E(T_E, q_2; 0, q_1) \approx \psi_0(q_2) \psi_0(q_1) \exp(-E_0 T_E); \quad (\text{when } T_E \rightarrow \infty) . \quad (18.4)$$

We now put $q_2 = q_1 = 0$, take the logarithm of both sides and divide by T_E ; then in the limit of $T_E \rightarrow \infty$, we get a formula for the ground state energy:

$$-E_0 = \lim_{T_E \rightarrow \infty} \left[\frac{1}{T_E} \ln K_E(T_E, 0; 0, 0) \right] . \quad (18.5)$$

Useful result:

1. Ground state energy

So, if we can determine the kernel by some method, we will know the ground state energy of the system. Once the ground state energy is known, we can plug it back into the asymptotic expansion in Eq. (18.4) and determine the ground state wavefunction.

Very often, we would have arranged matters such that the ground state energy of the system is actually zero. When $E_0 = 0$, there is a nice way of determining the wavefunction from the kernel by noting that:

$$\lim_{T \rightarrow \infty} K(T, 0; 0, q) \approx \psi_0(0) \psi_0(q) \propto \psi_0(q) . \quad (18.6)$$

So, the infinite time limit of the kernel — once we have introduced the imaginary time — allows determination of both the ground state wavefunction as well as the ground state energy. The proportionality constant ψ_0 can be fixed by normalizing the wavefunction.

Useful result:

2. Ground state wavefunction

Of course, these ideas are useful only if we can compute the kernel without knowing the wavefunctions in the first place. This is possible — as we discussed in Chapter 17 — whenever the action is quadratic in the dynamical variable. In that case, the kernel in real time can be expressed in the form

$$K(t_2, q_2; t_1, q_1) = N(t_1, t_2) \exp i \mathcal{A}_c(t_2, q_2; t_1, q_1) , \quad (18.7)$$

where \mathcal{A}_c is the action evaluated for a classical trajectory and $N(t_2, t_1)$ is a normalization factor and we are using units with $\hbar = 1$. The same ideas will work even when we can *approximate* the kernel by the above expression. We saw in Chapter 17 that in the semiclassical limit, the wavefunctions can be expressed in terms of the classical action. It follows that the kernel can be written in the above form in the same semiclassical limit. If

we now analytically continue this expression to imaginary values of time, then, using the result in Eq. (18.6) we get a simple formula for the ground state wavefunction in terms of the Euclidean action (that is, the action for a classical trajectory obtained after analytic continuation to imaginary values of time):

$$\psi_0(q) \propto \exp[-\mathcal{A}_E(T_E = \infty, 0; T_E = 0, q)] \propto \exp[-\mathcal{A}_E(\infty, 0; 0, q)] . \quad (18.8)$$

As an application of these results, consider a simple harmonic oscillator with the Lagrangian $L = (1/2)(\dot{q}^2 - \omega^2 q^2)$. The classical action with the boundary conditions $q(0) = q_i$ and $q(T) = q_f$ is given by

Example, you know what!

$$\mathcal{A}_c = \frac{\omega}{2 \sin \omega T} [(q_i^2 + q_f^2) \cos \omega T - 2q_i q_f] . \quad (18.9)$$

Analytic continuation will give the Euclidean action corresponding to $i\mathcal{A}_c$ to be $-\mathcal{A}_E$ where

$$\mathcal{A}_E = \frac{\omega}{2 \sinh \omega T} [(q_i^2 + q_f^2) \cosh \omega T - 2q_i q_f] . \quad (18.10)$$

Using this in Eq. (18.8), we find that the ground state wavefunction has the form

$$\psi_0(q) \propto \exp[-(\omega/2)q^2] , \quad (18.11)$$

which, of course, is the standard result. You can also obtain the ground state energy $(1/2)\hbar\omega$ by using Eq. (18.5). What is amazing, when you think about it, is that the Euclidean kernel in the limit of an *infinite time interval* has information about the *ground state* of quantum system. This is the first example in which imaginary time leads to a real result!

The analytic continuation to imaginary values of time also has close mathematical connections with the description of systems in a thermal bath. To see this, consider the mean value of some observable $\mathcal{O}(q)$ of a quantum mechanical system. If the system is in an energy eigenstate described by the wavefunction $\psi_n(q)$, then the expectation value of $\mathcal{O}(q)$ can be obtained by integrating $\mathcal{O}(q)|\psi_n(q)|^2$ over q . If the system is in a thermal bath at temperature β^{-1} , described by a canonical ensemble, then the mean value has to be computed by averaging over all the energy eigenstates *as well* with a weightage $\exp(-\beta E_n)$. In this case, the mean value can be expressed as

This one is unexpected

Thermal + Quantum average

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \sum_n \int dq \psi_n(q) \mathcal{O}(q) \psi_n^*(q) e^{-\beta E_n} \equiv \frac{1}{Z} \int dq \rho(q, q) \mathcal{O}(q) , \quad (18.12)$$

where Z is the partition function and we have defined a *density matrix* $\rho(q, q')$ by

$$\rho(q, q') \equiv \sum_n \psi_n(q) \psi_n^*(q') e^{-\beta E_n}, \quad (18.13)$$

in terms of which we can rewrite Eq. (18.12) as

$$\langle \mathcal{O} \rangle = \frac{\text{Tr}(\rho \mathcal{O})}{\text{Tr}(\rho)}, \quad (18.14)$$

where the trace operation involves setting $q = q'$ and integrating over q .

This standard result shows how $\rho(q, q')$ contains information about both thermal and quantum mechanical averaging. In fact, the expression for the density matrix in Eq. (18.13) is the coordinate basis representation of the matrix corresponding to the operator $\rho = \exp(-\beta H)$. That is,

$$\rho(q, q') = \langle q | e^{-\beta H} | q' \rangle. \quad (18.15)$$

But what is interesting is that we can now relate the density matrix of a system in finite temperature — something very real and physical — to the path integral kernel in imaginary time! This is obvious from comparing Eq. (18.13) with Eq. (18.1). We find that the density matrix can be immediately obtained from the Euclidean kernel by:

$$\rho(q, q') = K_E(\beta, q; 0, q'). \quad (18.16)$$

Temperature from imaginary time!

The imaginary time is now being identified with the inverse temperature. Very crudely, this identification arises from the fact that thermodynamics in the canonical ensemble uses $e^{-\beta H}$ while the standard time evolution in quantum mechanics uses e^{-iH} . But beyond that, it is difficult to understand in purely physical terms why imaginary time and real temperature should have anything to do with each other.

In obtaining the expectation values of operators which depend only on q — like the ones used in Eq. (18.12) — we only need to know the diagonal elements $\rho(q, q) = K_E(\beta, q; 0, q)$. The kernel in the right hand side can be thought of as the one corresponding to a periodic motion in which a particle starts and ends at q in a time interval β . *In other words, periodicity in imaginary time is now linked to finite temperature.*

Bonus: Black hole temperature in just two steps!

Believe it or not, most of the results in black hole thermodynamics can be obtained *from this single fact* by noting that the spacetimes representing a black hole, for example, have the appropriate periodicity in imaginary time. Considering the elegance of this result, let us pause for a moment and see how it comes about. Consider a curved spacetime in general relativity which has a line interval

$$ds^2 = -f(r)dt^2 + \frac{dr^2}{f(r)} + dL_\perp^2, \quad (18.17)$$

where dL_{\perp}^2 represents the metric in the two transverse directions. For example, we saw in Chapter 11 that the Schwarzschild metric representing a black hole has this form with $f(r) = 1 - (r_g/r)$ where $r_g = (2GM/c^2) = 2M$ (in units with $G = c = 1$) and dL_{\perp}^2 represents the standard metric on a two-sphere. The only property we will actually need is that $f(r)$ has a simple zero at some $r = a$ with $f'(a) \equiv 2\kappa$ being some constant. In the case of the black hole metric, $\kappa = (1/2r_g)$. When we consider the metric near the horizon $r \approx a$, we can expand $f(r)$ in a Taylor series and reduce it to the form

$$ds^2 = -2\kappa l dt^2 + \frac{dl^2}{2\kappa l} + dL_{\perp}^2, \quad (18.18)$$

Step 1: Metric near a horizon

where $l \equiv (r - a)$ is the distance from the horizon. If we now make a coordinate transformation from l to another spatial coordinate x such that $(\kappa x)^2 = 2\kappa l$, the metric becomes

$$ds^2 = -\kappa^2 x^2 dt^2 + dx^2 + dL_{\perp}^2. \quad (18.19)$$

This represents the metric near the horizon of a black hole.

So far we have not done anything non-trivial. Now we shall analytically continue to imaginary values of time with $it = \tau$ and denote $\kappa\tau = \theta$. Then the corresponding analytically continued metric becomes

Step 2: Go Euclidean; find the period of imaginary time

$$ds^2 = x^2 d\theta^2 + dx^2 + dL_{\perp}^2. \quad (18.20)$$

But $(dx^2 + x^2 d\theta^2)$ is just the metric on a two dimensional plane in polar coordinates and if it has to be well behaved at $x = 0$, the coordinate θ must be periodic with period 2π . Since $\theta = \kappa\tau$, it follows that the imaginary time τ must be periodic with period $2\pi/\kappa$ as far as any physical phenomenon is concerned. But we saw earlier that such a periodicity of the imaginary time is mathematically identical to working with finite temperature, with the temperature

$$\beta^{-1} = \frac{\kappa}{2\pi} = \frac{1}{4\pi r_g} = \frac{\hbar c^3}{8\pi GM}, \quad (18.21)$$

where the first equality is valid for a general class of metrics (with κ suitably defined by Taylor expansion of $f(r)$) while the last two results are for the Schwarzschild metric, and in the final expression, we have reverted back to normal units. This is precisely the Hawking temperature of a black hole of mass M which we obtained by a different method in Chapter 12. Here we could do that just by looking at the form of the metric near the horizon and using the relation between periodicity in imaginary time and temperature.

There you are!

The imaginary time and Euclidean action also play an interesting role in the case of tunneling. To see this, let us start with the expression for the

Another application classical action written in the Jacobi-Mapertius form (see Eq. (2.37):

$$S = \int p dq = \int \sqrt{2m(E - V)} dq . \quad (18.22)$$

As long as $E > V$, this will lead to a real value for S . Tunneling occurs, however, when $E < V$. To simplify matters a little bit, let us consider the case of a particle with $E = 0$ (which can always be achieved by adding a constant to the Hamiltonian) moving in a potential $V > 0$. In that case the action becomes pure imaginary and is given by

$$S = i \int \sqrt{2mV} dq , \quad (18.23)$$

and the corresponding branch of the semiclassical wavefunction will be exponentially damped:

$$\psi \propto \exp iS = \exp - \int \sqrt{2mV} dq . \quad (18.24)$$

This represents the fact that you cannot have a classical trajectory with $E = 0$ in a region in which $V > 0$.

*What you can't do
in real time, you can
do in imaginary time*

It is however possible to have such a trajectory if we analytically continue to imaginary values of time. In real time, the conservation of energy for a particle with $E = 0$ gives $(1/2)m(dq/dt)^2 = -V(q)$ which cannot have real solutions when $V > 0$. But when we set $t = -i\tau$ this equation becomes $(1/2)m(dq/d\tau)^2 = V(q)$ which, of course, has perfectly valid solutions when $V > 0$. So the tunneling through a potential barrier can be interpreted as a particle moving off to imaginary values of time as far as the mathematics goes. The *Euclidean* action will now be

$$S_E = \int \sqrt{2mV} dq . \quad (18.25)$$

All that we need to do to obtain the tunneling amplitude is to replace iS by $-S_E$ in the argument of the relevant exponential so that the wavefunction in Eq. (18.24) becomes:

$$\psi \propto \exp iS = \exp - \int \sqrt{2mV} dq = \exp -S_E . \quad (18.26)$$

So we find that the tunneling amplitude across the potential can also be related to analytic continuation in the imaginary time and and the Euclidean action.

Schwinger effect

We will now use these ideas to obtain a really non-trivial phenomenon in quantum electrodynamics, called the Schwinger effect, named after Julian Schwinger who was one of the creators of quantum electrodynamics and received a Nobel Prize for the same. In simplest terms, this effect

can be stated as follows. Consider a region of space in which there exists a constant, uniform electric field. One way to do this is to set-up two large, parallel, conducting plates separated by some distance L and connect them to the opposite poles of a battery. This charges the plates and produces a constant electric field between them. Schwinger showed that, in such a configuration, electrons and positrons will spontaneously appear in the region between the plates through a process which is called pair production from the vacuum.

The first question one would ask is how particles can appear out of nowhere. This is natural since we haven't seen tennis balls or chairs appear out of the vacuum spontaneously. In quantum field theory, what we call vacuum is actually bristling with quantum fluctuations of the fields which can be interpreted in terms of virtual particle-antiparticle pairs. Under normal circumstances, such a virtual electron-positron pair will be described by the situation in the left frame of Fig. 18.1. We think of an electron and positron being created at the event A and then getting annihilated at the event B . In the absence of any external fields, there is no force acting on these virtual pairs and they continuously appear and disappear quite randomly in the spacetime.

How can they pop out of the vacuum?

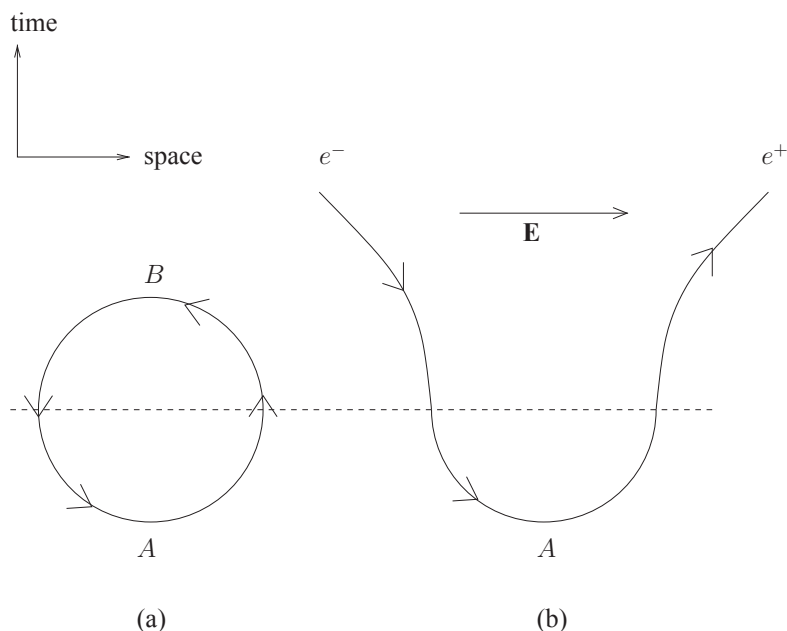


Fig. 18.1: In the vacuum, there exist virtual electron-positron pairs which are constantly created and annihilated as shown in the left frame (a). An electron-positron pair is created at A and annihilated at B with the positron being interpreted as an electron going backward in time. The right frame (b) shows how, in the presence of an electric field, this virtual process can lead to creation of *real* electrons and positrons.

One way to think about it

Consider now what happens if there is an electric field present in this region of space. The electric field will pull the electron in one direction and push the positron in the opposite direction since the electrons and positrons carry opposite charges. In the process, the electric field will do work on the virtual particle-antiparticle pair and hence will supply energy to them. If the field is strong enough, it can supply an energy greater than the rest energy of the two charged particles which is just $2 \times mc^2$ where m is the mass of the particle. This allows the virtual particles to become real. That is how the constant electric field between two conducting parallel plates produces particles out of the vacuum. It essentially does work on the virtual electron-positron pairs which are present in the space-time and converts them into real particles as shown in the right frame of Fig. 18.1(b).

One way to model this is to assume that the particle tunnels from the trajectory on the left to the one on the right through the semicircular path in the lower half. The trajectories on the left and right are real trajectories for the charged particle but the semicircle is a ‘forbidden’ quantum process. We will now see how the imaginary time makes this possible.

Imaginary time makes virtual, real

To do this, we begin with the trajectory in real time which will correspond to relativistic motion with uniform acceleration $g = qE/m$. We have worked this out in Chapter 12 and the result is given — with suitable choice of initial conditions — by:

$$x = (1/g) \cosh(g\tau); \quad t = (1/g) \sinh(g\tau); \quad x^2 - t^2 = 1/g^2. \quad (18.27)$$

The trajectory is a (pair of) hyperbola in the $t - x$ plane shown in Fig. 18.1(b). If we now analytically continue to imaginary values of τ and t , the trajectory becomes a circle $x^2 + t_E^2 = 1/g^2$ of radius $(1/g)$ and the parametric equations become

$$x = (1/g) \cos \theta; \quad t = (1/g) \sin \theta; \quad \theta = g\tau_E. \quad (18.28)$$

By going from $\theta = \pi$ to $\theta = 0$, say, we can get this to be a semicircle connecting the two hyperbolas.

Back to action

To obtain the amplitude for this process we have to evaluate the value of the Euclidean action for the semicircular track. The action for a particle of charge q in a constant electric field E represented by a scalar potential $\phi = -Ex$ is given by

$$A = -m \int d\tau + qE \int x dt, \quad (18.29)$$

where τ is the proper time of the particle. So, on analytic continuation we get

$$iA = -im \int d\tau + iqE \int x dt \rightarrow -m \int d\tau_E + qE \int x dt_E \equiv -A_E. \quad (18.30)$$

The Euclidean action A_E in Eq. (18.30) can be easily transformed to an integral over θ and noting that the integral over $x dt_E$ is essentially the area enclosed by the curve, which is a semicircle of radius $(1/g)$, we get

$$-A_E = -\frac{m}{g} \int_{\pi}^{2\pi} d\theta + \frac{m}{2g} \int_{\pi}^{2\pi} d\theta = -\frac{m\pi}{2g} . \quad (18.31)$$

The limits of the integration are so chosen that the path in the imaginary time connects $x = -(1/g)$ with $x = (1/g)$ thereby allowing a virtual semi-circular loop to be formed as shown in Fig. 18.1(b). Hence the final result for the Euclidean action for this classically forbidden process is

$$A_E = \frac{\pi m}{2g} = \frac{\pi m^2}{2qE} . \quad (18.32)$$

With the usual rule that a process with $\exp iA$ gets replaced by $\exp(-A_E)$ when it is classically forbidden, we find the amplitude for this process to take place to be $\mathcal{A} \propto \exp(-A_E)$. The corresponding probability $\mathcal{P} = |\mathcal{A}|^2$ is given by

$$\mathcal{P} \approx \exp -(\pi m^2/qE) . \quad (18.33)$$

This is the leading term for the probability which Schwinger obtained for the pair creation process. (In fact, one can even obtain the sub-leading terms by transferring paths which wind around several times in the circle but we will not go into this; if you are interested, take a look at ref. [70]). Once again, the moral is clear. What is forbidden in real time is allowed in imaginary time!

Final result

The expression for \mathcal{P} is non-analytic in q which measures the strength of coupling between the charge and the electromagnetic field. Usually, in quantum field theory one studies processes (like e.g., scattering) by a perturbative expansion in q . It is obvious that you will not be able to calculate \mathcal{P} by such a procedure, irrespective of how many orders in perturbation you calculate! The approach based on Euclidean action is capable of giving us non-perturbative results.

You can't get it in perturbation theory!

So far we were making things complex by analytically continuing from real to imaginary *time*. There are other physical situations in which this idea does not work but you can get around by actually using a complex *coordinate* (rather than time). One beautiful application of this technique is in understanding a phenomenon called over-the-barrier-reflection in a potential. Let me describe this situation which, somehow, does not find adequate discussion in textbooks.

Complex space?

Consider a potential of the form in Fig. 18.2 in which a particle is incident from the left. If its energy is like E_0 , which is below the peak of the potential, it will tunnel to the right and we have already seen that one can obtain the transmission coefficient T by analytically continuing

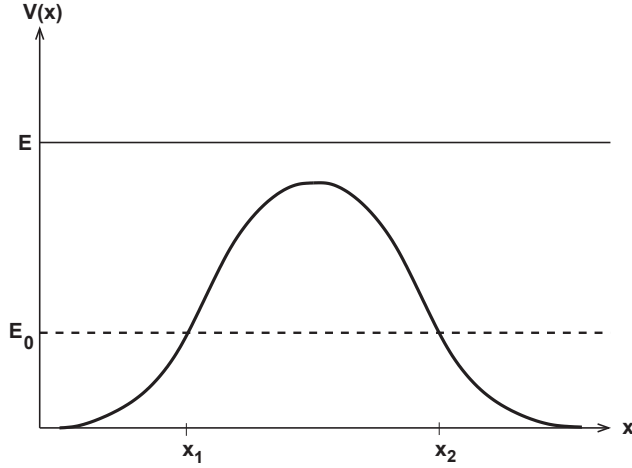


Fig. 18.2: A generic potential indicating the energy levels at which (i) tunneling and (ii) over-the-barrier-reflection can occur. (i) A particle incident from the left with energy E_0 will tunnel through the potential with an exponentially small transmission coefficient; its reflection coefficient will be nearly unity. (ii) On the other hand, a particle incident from the left with energy E will be *reflected* with an exponentially small amplitude; its transmission coefficient will be nearly unity.

to complex *time*. This T is exponentially small and is non-analytic in \hbar ; it goes to the credit of imaginary time method that we can pick it up. Of course, the standard WKB method — which involves the integral of $p(E_0, q)dq$ — will also lead to the correct result in this case because p will become imaginary when $E_0 < V$.

Over-the-barrier-reflection, described

Consider next a particle with energy E (as shown in Fig. 18.2) which is flying above the peak of the potential. Classically, the transmission coefficient is now unity and the reflection coefficient is zero. But quantum mechanically, we know that there is a small reflection coefficient $R \neq 0$ which is now exponentially small. As an example, consider a potential (chosen because the exact solution is known!) of the form

$$V(x) = \frac{V_0}{1 + e^{-x/a}}. \quad (18.34)$$

The reflection coefficient for this case happens to be

$$R = \frac{\sinh^2 \pi(k_1 - k_2)a}{\sinh^2 \pi(k_1 + k_2)a}; \quad T = 1 - R, \quad (18.35)$$

where we have defined

$$k_1 = \frac{1}{\hbar} \sqrt{2mE}; \quad k_2 = \frac{1}{\hbar} \sqrt{2m(E - V_0)}. \quad (18.36)$$

How do we get this result from a WKB like approximation? The $p(E, q)$ remains real now and hence integrating $p dq$ over any range of real q will not lead to tunnelling probability. So it is obvious that going to imaginary time is not going to work and a different trick is needed. What we need to do, is to go to complex coordinates and look at the paths in the complex plane [70, 71].

To illustrate this procedure, it will be useful to consider the turning points for the problem defined by the equation $E = V(z)$ where we have now analytically continued from real x to complex z . When the energy is like E_0 in Fig. 18.2, we see that there are two turning points, both of which are real indicated by x_1 and x_2 in Fig. 18.2. What is more, there is a branch cut in the real axis in the complex plane between x_1 and x_2 . The standard tunneling problem now corresponds to integrating through the potential along the path C_1 shown in Fig. 18.3. You can convince yourself that this will give the correct result.

What happens in the complex plane

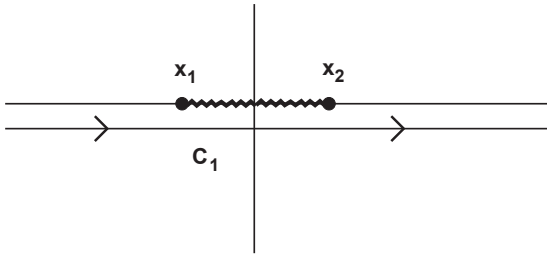


Fig. 18.3: Tunneling through a potential by a particle with energy E_0 in Fig. 18.2 can be described using the contour C_1 . The turning points x_1 and x_2 (where $E = V(x)$) are on the real axis with a branch cut connecting them in the complex plane.

As we increase the energy from E_0 , the turning points approach each other and coalesce together at some point when the energy is just equal to the maximum of the potential. When the energy increases further so that all regions are classically accessible, there are no *real* turning points. The equation $E = V(z)$ will, of course, have complex solutions. We will pick the complex solution for which the turning point is closest to the real axis. For illustration, consider a situation like the one shown in Fig. 18.4. The branch cuts are now on the imaginary axis for the simplest case one can consider.

We want a rule to determine the exponentially small reflection coefficient in this particular case. This rule is essentially based on distorting the path in the complex plane in the form of the curve C_2 shown in Fig. 18.4. The reflection coefficient is now given by the expression

Rule for getting reflection coefficient

$$R = \left| \exp \left(\int_{C_2} k(z) dz \right) \right|^2. \quad (18.37)$$

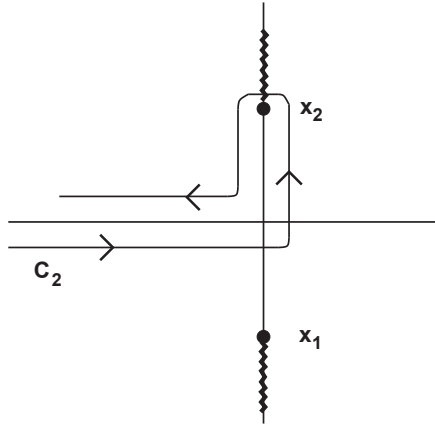


Fig. 18.4: The over-the-barrier-reflection by a particle with energy E in Fig. 18.2 can be described using the contour C_2 . The turning points x_1 and x_2 (where $E = V(x)$) are now on the *imaginary* axis with the branch cuts as shown in the figure.

Let us illustrate this for the case of the potential in Eq. (18.34) for which the branch points, in the case of $E > V_0$, occur at

$$x_c = -a \ln \left(1 - \frac{V_0}{E} \right) \pm ia(2n+1)\pi; \quad n = 0, 1, 2, \dots \quad (18.38)$$

The contribution from the branch point closest to the real axis (viz., the one with $n = 0$) will dominate the result and the rest will be exponentially small and can be ignored. The choice of the contour in Fig. 18.4 shows that the path is ascending on the first sheet and descending in the second which ensures that $R < 1$. Then the relevant WKB integral along the contour has the form

$$\int_{C_2} k(z) dz = \phi_1 + 2i\sigma_1, \quad (18.39)$$

where ϕ_1 is real and

$$\sigma_1 = ik_1 \int_{x_0}^{x_0 + i\pi a} dx \left(1 - \frac{1}{E} V(x) \right)^{1/2}, \quad (18.40)$$

with

$$x_0 = \text{Re } x_c = -a \ln \left(1 - \frac{V_0}{E} \right). \quad (18.41)$$

It is clear from our expression for the reflection coefficient Eq. (18.37), that ϕ_1 does not contribute, and, in σ_1 , only the real part makes a contri-

bution. With some tricks in contour integration, we can easily show that

*Some maths,
slightly tricky*

$$\operatorname{Re} \sigma_1 = \pi k_1 a \left(1 - \frac{V_0}{E}\right)^{1/2} = \pi k_2 a . \quad (18.42)$$

One way to proceed is as follows. The contour integral we needed to calculate is of the form

$$I = \int_{x_0}^{x_0 + i\pi a} dx \left(1 - \frac{\rho}{1 + e^{-x/a}}\right)^{1/2}, \quad (18.43)$$

where

$$x_0 = -a \ln \left(1 - \frac{V_0}{E}\right), \quad 0 < \rho = \frac{V_0}{E} < 1 . \quad (18.44)$$

To evaluate this, we introduce the variable

$$z = 1 - \frac{\rho}{1 + e^{-x/a}}, \quad (18.45)$$

which converts the integral to the form

$$I = a \int_0^{2\left(\frac{1-\rho}{2-\rho}\right)} dz \left(\frac{1}{1-z} + \frac{1}{\rho-1+z}\right). \quad (18.46)$$

The singularity is now at $z = (1 - \rho)$. If we deform the contour around this point by a tiny semicircle, we pick up the imaginary contribution:

$$\operatorname{Im} I = \pi a \sqrt{1 - \rho}, \quad (18.47)$$

so that

$$\operatorname{Re} \sigma_1 = \operatorname{Im} k_1 I = \pi a k_2 . \quad (18.48)$$

Substituting this result in Eq. (18.37), we get our reflection coefficient:

Final result

$$R = \left| \exp \left(i \int_{C_4} k(x) dx \right) \right|^2 = e^{-4 \operatorname{Re} \sigma_1} = e^{-4\pi k_2 a} . \quad (18.49)$$

As a check, we can compare it with the leading term of the exact result in Eq. (18.35) which is given by

$$\lim_{\hbar \rightarrow 0} \frac{\sinh^2 \pi(k_1 - k_2)a}{\sinh^2 \pi(k_1 + k_2)a} = \frac{\exp 2\pi(k_1 - k_2)a}{\exp 2\pi(k_1 + k_2)a} = \exp(-4\pi k_2 a) . \quad (18.50)$$

Clearly, we got the correct result.

Verified: all is fine

It is interesting that this procedure involving the complex path picks out the exponentially small reflection coefficient which is non-analytic in \hbar . In fact, this procedure of the complex path is more general than you might think. By choosing the contours appropriately, we can also directly obtain the transmission coefficient in this particular case (which, of course is nearly unity because $T = (1 - R)$). You can also work out the transmission and reflection coefficient in the usual tunneling case by the complex path method. But in this case, the complex *time* procedure is much more natural.

An electromagnetic field can exert a force on a charged particle. The expression for this force, viz. the Lorentz force, is given by $q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$. This tells you that, if \mathbf{E} and \mathbf{B} vanish, then there is no force. Similarly, we do not expect any Lorentz force to arise when there are no charged particles. All these sound quite reasonable.

*Classical physics:
If it is not there, it
is not there*

Unfortunately, the real world is quantum mechanical, governed by quantum amplitudes, the uncertainty principle and what not. The quantum theory of the electromagnetic field leads to the notion of photons, and the closest description of the “absence of electromagnetic field” would correspond to a quantum state with zero photons, which is usually called the vacuum state of the electromagnetic field. One would have imagined that, if there are no photons, then there will be no measurable physical effects due to the electromagnetic field. While this is *more or less* true — which is rather reassuring — there are indeed interesting situations in which it is *not* true! We will describe [72] one such context, called the Casimir effect, in this chapter.

*Quantum Physics:
If it is not there, it
may still be there*

The simplest — though a bit idealized — description of the Casimir effect is the following. Consider two parallel, perfectly conducting, plates kept in otherwise empty space, separated by a distance L . Then, they will attract each other with a force

An eerie force

$$F = -\frac{\pi^2}{240} \frac{\hbar c}{L^4} A, \quad (19.1)$$

where A is the cross-sectional area of either plate!! Note that there are no net charges put on the plates; we are *not* talking about a *charged* parallel plate capacitor. The force acts between two plates kept in the vacuum. This effect was predicted [73] by the Dutch physicist Hendrick Casimir in 1948 and has been measured in the lab [74, 75]. One nice way of understanding this result is in terms of a tangible force exerted by the electromagnetic *vacuum*. Let us see how.

One way to see it coming

Before launching into mathematics, let us try to understand the basic reason for this phenomenon in qualitative terms. Consider the familiar example of a harmonic oscillator, with the Hamiltonian

$$H(p, q) = \frac{1}{2}[p^2 + \omega^2 q^2]. \quad (19.2)$$

We have set the mass of the particle to unity for simplicity. Classically, the minimum energy for such a system is zero ($E_{\text{class}} = 0$), which is achieved when $q = p = 0$. We know, however, that this is not possible in quantum theory, essentially because of uncertainty principle. To minimize the potential energy, we need to set $q = 0$; but if we know the position to such infinite precision, the momentum will be infinitely uncertain and we cannot guarantee a low value for $p^2/2$! So to minimize the total energy, we need to allow for some amount of fluctuation in both q and p that is commensurate with uncertainty principle $\Delta p \Delta q \gtrsim \hbar$. The resulting ground state will then have a non-zero energy $E_{\text{quant}} = (1/2)\hbar\omega$.

Suppose we consider a different physical system with the Hamiltonian $H_{\text{new}} = H(p, q) - (1/2)\hbar\omega$ where $H(p, q)$ is given by Eq. (19.2). Since the subtraction of a constant from the Hamiltonian does not change the equations of motion, we still again have a harmonic oscillator but with a shift in the energy. Classically, the minimum energy state will still correspond to $q = p = 0$ but with energy $E_{\text{class}} = -(1/2)\hbar\omega$. However, quantum mechanically, the ground state will exhibit fluctuations in q and p but this state will now have zero energy; $E_{\text{quant}} = 0$! This is the crucial point. Quantum mechanics allows you to have a state for the harmonic oscillator with the Hamiltonian $H_{\text{new}}(p, q)$ such that $E_{\text{quant}} = 0$, which *can* host fluctuations in the dynamical variables q and p .

The q and p can fluctuate even if $E = 0$; don't confuse these two!

Something very analogous happens in the case of an electromagnetic field. As we shall see, the electromagnetic field can be thought of as a bunch of harmonic oscillators. The ground state will correspond to a state of zero photons and one can also arrange matters such that it has zero energy. But the electric and magnetic fields will play roles analogous to p and q of the oscillator and they will exhibit fluctuations in the ground state — which are usually called vacuum fluctuations. Therefore, one cannot really say that the electromagnetic fields vanish in the vacuum state even though we can make its energy vanish. This is completely analogous to the fact that we cannot say the position q of the oscillator vanishes in the ground state of the oscillator. Once we recognize this fact, it is not surprising that the electromagnetic vacuum can exert forces on bodies. In real life, the situation is a little bit more complicated because the procedure analogous to the subtraction of $(1/2)\hbar\omega$ is more non-trivial in the case of the electromagnetic field; but the essential idea is the same.

Why vacuum fluctuations arise

Two simplifications

Let us now try to understand this mathematically. As it turns out, the essential idea can be illustrated by ignoring two complications of the real

world. The first is the the vector nature of electromagnetism and the second is the fact that space is three dimensional. We will first work out a simpler picture using a scalar field with just one degree of freedom (rather than with the electromagnetic field) and also ignoring the two transverse directions and treating space as one-dimensional. After we work out the simplified picture, we will describe how it generalizes to the real case.

Once we ignore the vector nature of the electromagnetic field, we can work with a single scalar field $\phi(t, x)$ which — if you want — can be thought of as analogous to any one component of the electromagnetic field. In the absence of sources, we know that each component of the field satisfies the wave equation, which can now be written as:

$$\frac{\partial^2 \phi}{\partial t^2} - \frac{\partial^2 \phi}{\partial x^2} = 0 . \quad (19.3)$$

In three dimensions, the second term would have been $-\nabla^2 \phi$ which becomes one-dimensional when we ignore two spatial coordinates. We have also chosen units with $c = 1$. This equation can be simplified by introducing the spatial Fourier transform Q_k of $\phi(t, x)$ by

$$\phi(t, x) = \int_{-\infty}^{\infty} \frac{dk}{(2\pi)} Q_k(t) \exp[ikx] . \quad (19.4)$$

Substituting this in Eq. (19.3), we find that $Q_k(t)$ satisfies the equation $\ddot{Q}_k + k^2 Q_k = 0$. The field $\phi(t, x)$ is completely specified by the function $Q_k(t)$ so that we can think of $Q_k(t)$ as the dynamical variables describing our system. The fact that we are dealing with the field translates into the fact that we now have an *infinite* number of dynamical variables, one for each value of k . Other than that, we can work directly with $Q_k(t)$ instead of the original field $\phi(t, x)$.

*Degrees of freedom
in Fourier space*

One minor problem with $Q_k(t)$ is that it is a complex number (since $\phi(t, x)$ is real) and we would like to work with dynamical variables that are real. This is easily taken care of. As Q_k is complex, we have two degrees of freedom corresponding to the real and imaginary parts of Q_k for each k with the constraint $Q_k^* = Q_{-k}$. If we write $Q_k = (A_k + iB_k)$, then, since ϕ is a real scalar field, we can relate the variables for k to that for $-k$ as $A_k = A_{-k}$ and $B_k = -B_{-k}$. Evidently, only half the modes constitute independent degrees of freedom. Therefore, we can work with a new set of real modes q_k , defined for all values of k with a suitable redefinition, say, by taking $q_k = A_k$ for one half of k and $q_{-k} = B_k$ for the other half. This will, of course, lead to the same equation but for the *real* variable $q_k(t)$:

*Minor irritation,
disposed of*

$$\ddot{q}_k + k^2 q_k = 0 . \quad (19.5)$$

That is, the dynamical variable $q_k(t)$ satisfies the harmonic oscillator equation with frequency $\omega = |k|$, for each value of k . Our field is math-

If you can do oscillators, you can do field theory

ematically the same as an infinite number of harmonic oscillators, one for each k . It follows that everything we know about harmonic oscillators can now be applied to this system. In particular, we can quantize the field by quantizing each of the harmonic oscillators $q_k(t)$. (In fact, that is the essence of quantum field theory of non-interacting fields; the rest is just detail.)

Ground state of the oscillator ...

Classically, we can now construct the ground state by taking $q_k = 0$ for all values of k . This will, of course, make the field vanish along with its energy, as is to be expected from a sensible ground state. But, as we discussed earlier, this does not hold for the quantum ground state. The ground state of the harmonic oscillator for a given value of k is described by the ground state energy eigenfunction

$$\psi(q_k) = \left(\frac{\omega_k}{\pi}\right)^{1/4} \exp\left(-\frac{1}{2}\omega_k q_k^2\right). \quad (19.6)$$

The ground state wavefunction for the full system, made of a bunch of independent oscillators, can be described by the product of the ground state wavefunctions of each of the oscillators:

$$\Psi[\phi(x)] = \prod_k \left(\frac{\omega_k}{\pi}\right)^{1/4} \exp\left(-\frac{1}{2}\omega_k q_k^2\right) \equiv \bar{N} \exp\left[-\frac{1}{2} \int_{-\infty}^{\infty} \frac{dk}{(2\pi)} \omega_k q_k^2\right]. \quad (19.7)$$

... gives you the ground state of the field theory

This expression can be interpreted along similar lines as the harmonic oscillator wavefunction in usual quantum mechanics. Suppose we have a harmonic oscillator in the ground state and we measure the position q . Then the relative probability that we will get a value $q = a$, compared to a value $q = b$ is given by:

$$R = \frac{|\psi(a)|^2}{|\psi(b)|^2} = \exp(-\omega[a^2 - b^2]). \quad (19.8)$$

Now suppose we have a *quantum field* which is in the ground state and we measure the field everywhere at, say, $t = 0$. Then, there is some probability that we will get a field configuration described by the function $\phi(0, x) = f_1(x)$ and some other probability that field configuration is described by the function $\phi(0, x) = f_2(x)$. Just as in the previous case, we want to know the relative probability of getting one configuration compared to another. To find this, we first obtain the spatial Fourier transforms of $f_1(x)$ and $f_2(x)$ and call them a_k and b_k . Then, the relative probability is given by

$$R = \frac{|\Psi(f_1(x))|^2}{|\Psi(f_2(x))|^2} = \exp\left(-\int_{-\infty}^{\infty} \frac{dk}{(2\pi)} \omega_k [|a_k|^2 - |b_k|^2]\right). \quad (19.9)$$

For any choice of $f_1(x)$ and $f_2(x)$ the above number can be computed, allowing us to determine the relative probability for the vacuum state to host two different field configurations.

You would have noticed that we switched to relative probabilities from absolute probabilities in this discussion. For a single harmonic oscillator, one could have said that $|\psi(q)|^2 dq$ gives the absolute probability of finding the particle in the interval $(q, q + dq)$. When we have an infinite number of oscillators, the normalization factor \bar{N} in Eq. (19.7) involves an infinite product which is hard to define rigorously. We bypass this by using relative probabilities, in which the normalization factor cancels out.

Why relative probability

Before we proceed further, let me mention the corresponding result in three spatial dimensions. In this case Eq. (19.7) has the obvious generalization to:

$$\Psi[\phi(\mathbf{x})] = \prod_k \left(\frac{\omega_k}{\pi} \right)^{1/4} \exp\left(-\frac{1}{2} \omega_k |q_k|^2\right) = \bar{N} \exp\left[-\frac{1}{2} \int \frac{d^3 \mathbf{k}}{(2\pi)^3} \omega_k |q_k|^2\right]. \quad (19.10)$$

In fact, in this case it is nicer to exhibit the result in terms of the field configuration itself by using $\omega_k = |\mathbf{k}|$ and $\omega_k |q_k|^2 = k^2 |q_k|^2 / |\mathbf{k}|$. Since $i\mathbf{k}q_k$ is essentially the Fourier spatial transform of $\nabla\phi$, we can easily obtain

$$\int \frac{d^3 \mathbf{k}}{(2\pi)^3} \omega_k |q_k|^2 = \int \frac{d^3 \mathbf{k}}{(2\pi)^3} \frac{|\mathbf{k}|^2 |q_k|^2}{|\mathbf{k}|} = \frac{1}{2\pi^2} \int d^3 \mathbf{x} \int d^3 \mathbf{y} \left\{ \frac{\nabla_{\mathbf{x}} \phi \cdot \nabla_{\mathbf{y}} \phi}{|\mathbf{x} - \mathbf{y}|^2} \right\}. \quad (19.11)$$

Substituting this into Eq. (19.10) and taking the modulus, we get the probability distribution in the ground state to be:

A nice result

$$\mathcal{P}[\phi(\mathbf{x})] = |\Psi[\phi(\mathbf{x})]|^2 = N \exp \left\{ -\frac{1}{2\pi^2} \int \int d^3 \mathbf{x} d^3 \mathbf{y} \frac{\nabla_{\mathbf{x}} \phi \cdot \nabla_{\mathbf{y}} \phi}{|\mathbf{x} - \mathbf{y}|^2} \right\}, \quad (19.12)$$

with $N = |\bar{N}|^2$. Once again, this expression shows clearly that the vacuum state of the field can host — what is usually called — zero point fluctuations of the field variable ϕ . The probability that one detects a particular field configuration $\phi(\mathbf{x})$ when the field is in the vacuum state can be obtained by evaluating the value of \mathcal{P} for this particular functional form $\phi(\mathbf{x})$. The result is independent of time because of the stationarity of the vacuum state. Given the ambiguity in the overall normalization factor N , this probability should again be interpreted as a relative probability. That is, the ratio $(\mathcal{P}_1 / \mathcal{P}_2)$ will give the relative probability between two field configurations characterized by the functions $\phi_1(\mathbf{x})$ and $\phi_2(\mathbf{x})$.

Let us now ask what happens if we introduce two perfectly conducting parallel plates into the vacuum. The fact that the plates are perfectly conducting requires the electromagnetic field — for which our $\phi(t, \mathbf{x})$ is a proxy — to satisfy some non-trivial boundary conditions at $x = 0$ and $x = L$ where the plates are located. For the scalar field, we can

Back to the Casimir effect

take the boundary condition to be that the field vanishes at the plates: $\phi(t, 0) = \phi(t, L) = 0$ in one spatial dimension. You cannot describe a field satisfying such a boundary condition using the Fourier integral in Eq. (19.4) with k taking all possible values in $-\infty < k < \infty$. Instead, we must restrict it to a discrete — though infinite — set of values given by $k = n(\pi/L)$ and write:

$$\phi(t, x) = \sum_{n=1}^{\infty} q_n(t) \sin \left[n \frac{\pi x}{L} \right] , \quad (19.13)$$

so that the boundary conditions at $x = 0$ and $x = L$ are satisfied. We still have to deal with an infinite number of oscillators, but their frequencies are now given by $\omega_n = k_n = n(\pi/L)$.

If we now work out the corresponding ground state, it will be different from the one described by Eq. (19.7) because the integral over k will be now replaced by the sum over n . This is needed because, our boundary condition tells us that we have prevented the ground state from having non-zero probability for field configurations which do not vanish at the plates. The introduction of the plates — through the change in the boundary condition — has changed the ground state.

What about the energy of the ground state with and without the plates? They will also be different. In the absence of plates, each harmonic oscillator contributes an energy $(1/2)\hbar\omega = (1/2)\hbar|k|$ leading to a total ground state energy per unit length of space to be the integral over all k of $(1/2)\hbar|k|$; so the energy in a region of length L will be:

$$E_0 = \frac{L}{(2\pi)} \int_{-\infty}^{\infty} dk \frac{1}{2} \hbar |k| = \frac{L}{(2\pi)} \int_0^{\infty} dk \hbar k . \quad (19.14)$$

This is manifestly infinite, essentially because there are an infinite number of harmonic oscillators. What about the ground state energy in the presence of the plates? This is given by the sum

$$E'_0 = \frac{1}{2} \sum_{n=0}^{\infty} \hbar \omega_n = \frac{1}{2} \sum_{n=0}^{\infty} \hbar (n\pi/L) , \quad (19.15)$$

which is also infinite, essentially given by the sum of all positive integers.

These infinities are bad news but there is a trick to get around them. As we said before, the equation of motion for the k -th oscillator will not change if we subtract $(1/2)\hbar\omega_k$ from the Hamiltonian, but it will “regularize” the ground state energy to zero. This is equivalent to looking at the difference $(E'_0 - E_0)$ as the physically relevant quantity. To study this, it is convenient to introduce in Eq. (19.14) a continuous variable n via the equation $k = (\pi/L)n$. Then, we get from Eq. (19.14) and Eq. (19.15):

$$(E'_0 - E_0) = \frac{\hbar\pi}{2L} \left[\sum_{n=0}^{\infty} n - \int_0^{\infty} dn n \right] . \quad (19.16)$$

The vacuum changes!

And so does its energy

An important trick; get used to it

You may think that this is not of much help because this is of the form $(\infty - \infty)$ which does not have a precise meaning. That is true but there are ways of giving meaning to such expressions in a fairly systematic manner. The simplest procedure is to consider, instead of the expression in Eq. (19.16), the expression:

*Illegal operations
made legal?*

$$(E'_0(\lambda) - E_0(\lambda)) \equiv \frac{\hbar\pi}{2L} \left[\sum_{n=0}^{\infty} n \exp(-n\lambda) - \int_0^{\infty} dn n \exp(-n\lambda) \right]. \quad (19.17)$$

Here we have multiplied both the expressions by a “regulator function” $\exp(-n\lambda)$ where λ is just a parameter. Both the expressions as well as their difference are now finite and the idea is to *first* compute the difference as a function of λ and then take the limit of $\lambda \rightarrow 0$ hoping for the best. That is, we *interpret* the expression in Eq. (19.16) as the limit of the expression in Eq. (19.17) when $\lambda \rightarrow 0$. I will let you work out the expressions. You should first get:

$$\sum_{n=0}^{\infty} n \exp(-n\lambda) = \frac{e^{-\lambda}}{(1 - e^{-\lambda})^2} = \frac{1}{\lambda^2} - \frac{1}{12} + \frac{\lambda^2}{240} + \mathcal{O}(\lambda^4), \quad (19.18)$$

which diverges when $\lambda \rightarrow 0$, as to be expected. Similarly,

$$\int_0^{\infty} dn n \exp(-n\lambda) = \frac{1}{\lambda^2}, \quad (19.19)$$

which also diverges when $\lambda \rightarrow 0$. But, surprisingly, the difference between Eq. (19.18) and Eq. (19.19) remains finite as $\lambda \rightarrow 0$:

$$\sum_{n=0}^{\infty} n \exp(-n\lambda) - \int_0^{\infty} dn n \exp(-n\lambda) = -\frac{1}{12} + \mathcal{O}(\lambda^2) \rightarrow -\frac{1}{12}, \quad (19.20)$$

when $\lambda \rightarrow 0$. This allows us to obtain the following remarkable result:

$$E(L) \equiv (E'_0 - E_0) \equiv \lim_{\lambda \rightarrow 0} (E'_0(\lambda) - E_0(\lambda)) = -\frac{\pi\hbar}{24L}. \quad (19.21)$$

So, we see that the ground state energy of the system with the plates — when regularized by subtracting away the energy in the absence of the plates — is a negative number and is inversely proportional to the separation between the plates! (You may wonder whether the result will change if you use a different cutoff function other than $\exp -n\lambda$ in Eq. (19.17). One can actually prove that it will not, as long as the function satisfies some reasonable conditions.)

*No, it is not in the
eye of the beholder.*

It is clear that the energy in Eq. (19.21) will lead to an attractive force $F = -(dE/dL) \propto L^{-2}$ between the plates, since reducing the separation

between the plates leads to the lowering of the energy. A more physical way of thinking about this result is as follows. If we change the separation between the plates by an amount ΔL , the energy of the configuration will change by $(dE/dL)\Delta L$ which must be the work done by the agency separating the plates. Equating it to $-F\Delta L$, where F is the force acting between the plates, we find that $F = -dE/dL$.

Back to real life

In the mythical world of one spatial dimension, the plates are zero-dimensional points which are not of practical use. The corresponding calculation for the electromagnetic field in 3-dimensions is more complicated algebraically but all the concepts remain the same. The final result in this case is an expression for energy per unit transverse area of the plates, given by:

$$\frac{(E'_0 - E_0)}{A} = -\frac{\pi^2 \hbar c}{720 L^3}, \quad (19.22)$$

where we have re-introduced the c factor. The force per unit area acting between the plates is now given by

$$\frac{F}{A} = -\frac{d}{dL} \frac{(E'_0 - E_0)}{A} = -\frac{\pi^2 \hbar c}{240 L^4}. \quad (19.23)$$

This tiny force has actually been measured in the lab!

Let me outline the steps involved in this derivation. For the electromagnetic field in (1+3) dimensions, we consider a region between two parallel conducting plates, each of area $L \times L$, separated by a distance a . We will assume that $L \gg a$ and will be interested in computing the force per unit area of the conducting plates by differentiating the corresponding expression for the zero-point energy with respect to a . As in the previous case, we want to compute the zero-point energy in the presence of the plates, and in their absence, in the 3-dimensional volume (L^2a), and compute the difference.

The energy contained in this region in the absence of the plates is given by the integral

$$\begin{aligned} E_0 &= 2 \int \frac{L^2 d^2 k}{(2\pi)^2} \int \frac{a dk_3}{2\pi} \left[\frac{1}{2} \sqrt{k_1^2 + k_2^2 + k_3^2} \right] \\ &= \int \frac{L^2 d^2 k}{(2\pi)^2} \int_0^\infty dn \sqrt{k_1^2 + k_2^2 + \left(\frac{n\pi}{a}\right)^2}, \end{aligned} \quad (19.24)$$

Vacuum without plates

where the overall factor 2 in front in the first line takes into account two polarizations and we have set $k_3 = (n\pi/a)$ with a continuum variable n to obtain the second line. (However, note that, due to the change of the limits of integration, $dk_3 = 2dn(\pi/a)$.) Writing the transverse component of the wave vector as $k_\perp^2 \equiv k_1^2 + k_2^2 \equiv (\pi/a)^2 \mu$, we can re-write this expression

as an integral over μ and n in the form

$$E_0 = \frac{L^2}{2a} \int_0^\infty k_\perp dk_\perp \int_0^\infty dn [\mu + n^2]^{1/2} = \frac{\pi^2 L^2}{4a^3} \int_0^\infty d\mu \int_0^\infty dn [\mu + n^2]^{1/2}. \quad (19.25)$$

The integrals in both Eq. (19.24) and Eq. (19.25) are, of course, are divergent — as to be expected.

Let us next consider the situation in the presence of conducting plates. This will require replacing the integral over n by a summation over n when $n \neq 0$. When $n = 0$, the corresponding result has to be multiplied by a factor $(1/2)$ because only one polarization state contributes. It is straightforward to determine the allowed modes in this case and you will find that there are two polarizations, each contributing the energy:

Vacuum with plates

$$\omega_{k,n} = \sqrt{k_1^2 + k_2^2 + \left(\frac{n\pi}{a}\right)^2}, \quad (19.26)$$

when $n \neq 0$, and one polarization contributing when $n = 0$. Therefore, the corresponding expression in the presence of the plates is given by

$$\frac{E'_0}{L^2} = \frac{\pi^2}{4a^3} \left[\sum_{n=1}^\infty \int_0^\infty d\mu (\mu + n^2)^{1/2} + \frac{1}{2} \int_0^\infty d\mu \mu^{1/2} \right]. \quad (19.27)$$

We can evaluate both the divergent integrals exactly as before using a cut-off function. But just for fun, we will do it in a different way which adds to the mystery!

The new idea is to give meaning to the integral in Eq. (19.27) directly, without subtracting anything. That is, you just compute the integral in the presence of plates and ‘regularize’ its divergence by a trick. To do this, consider the integral

$$\int_0^\infty d\mu (\mu + n^2)^{-\alpha} = \frac{1}{(\alpha - 1)} n^{-2(\alpha - 1)}, \quad (19.28)$$

which is well defined for sufficiently large α . We can therefore write

$$\int_0^\infty d\mu \mu^{-\alpha} = \lim_{\Lambda \rightarrow 0} \int_0^\infty d\mu (\mu + \Lambda^2)^{-\alpha} = \lim_{\Lambda \rightarrow 0} \left[\frac{1}{(\alpha - 1)} \Lambda^{-2(\alpha - 1)} \right]. \quad (19.29)$$

If you put $\alpha = -1/2$ in this relation, you can prove the incredible result:

Power of positive thinking

$$\int_0^\infty d\mu \mu^{1/2} = \lim_{\Lambda \rightarrow 0} \left[\frac{1}{(-3/2)} \Lambda^3 \right] = 0. \quad (19.30)$$

In the language of dimensional regularization, a pet trick in high energy physics, this means that several power law divergences can be “regular-

ized” to vanish. This means we need not worry about the second integral within the square bracket in Eq. (19.27). It follows that the quantity we need to evaluate in Eq. (19.27) is given by the expression

$$\sum_{n=1}^{\infty} \int_0^{\infty} d\mu (\mu + n^2)^{-\alpha} = \frac{1}{(\alpha-1)} \sum_{n=1}^{\infty} \frac{1}{n^{2(\alpha-1)}} = \frac{1}{(\alpha-1)} \zeta(2\alpha-2), \quad (19.31)$$

in the limit of $\alpha \rightarrow -(1/2)$ where we have introduced the Riemann zeta function,

$$\zeta(x) = \sum_{n=1}^{\infty} n^{-x}. \quad (19.32)$$

You can cheat very, very rigorously!

You will recognize that the expression in Eq. (19.31) is just $\zeta(-3)$ which is the sum of the cubes of all the integers! One can define this quantity by analytic continuation in the complex plane and then one obtains the result (see Appendix to learn this black magic)

$$\zeta(1-2k) = -\frac{B_{2k}}{2k}, \quad (19.33)$$

where B_k (called the k -th Bernoulli number) is defined through the series expansion:

$$\frac{t}{e^t - 1} = \sum_{k=0}^{\infty} \frac{t^k}{k!} B_k. \quad (19.34)$$

Using Eq. (19.31) and Eq. (19.33), we get the final result to be:

$$\frac{E'_0}{L^2} = \frac{\pi^2}{4a^3} \left[\frac{B_4}{6} \right] = \frac{\pi^2}{24a^3} \left(\frac{-1}{30} \right) = -\frac{\pi^2}{720a^3}, \quad (19.35)$$

where we have used the fact that $B_4 = -(1/30)$. The corresponding force of attraction (per unit area of the plates) is given by

$$f = -\frac{\partial \mathcal{E}}{\partial a} = -\frac{\pi^2}{240a^4}. \quad (19.36)$$

It is this force (which works out to 10^{-8} N for $a = 1 \mu\text{m}$, $L = 1$ cm) that has been observed in the lab.

Oh, these theorists!

Incidentally, with the same “regularization”, quantum field theorists often conclude that the sum of all positive integers is not only finite but is a negative fraction $(-1/12)$! The corresponding zeta function regularization will reproduce the $(1+1)$ scalar field result without, of course, making us any wiser as to what is going on.

Box 19.1: The Casimir Effect: Some History

The history of the Casimir effect is rather curious in some aspects. The effect gets its name from a paper published by H.B.T. Casimir in 1948 in the Proceedings of the Royal Academy of Sciences of Netherlands. This paper had the basic result that two metallic plates would attract each other, but it did not really draw much attention in the next two or three decades. In fact, another work by Casimir and Polder published in 1948 got a lot more attention, possibly because it was published in the Physical Review rather than a Dutch journal.

Another curiosity is that, in the early decades after the publication, the Casimir effect was considered neither spectacular nor mysterious. Recall that Casimir's paper came 75 years after the celebrated thesis of Van der Waals which introduced the weak attractive force between *neutral* molecules and 18 years after F. London gave an explanation for the Van der Waals force in terms of fluctuating electric dipole moments. So the fact that neutral metallic bodies — as long as they *internally* contain charged particles — could exert forces on each other, did not cause much surprise to the community.

Years later, in 1954, E.M. Lifshitz provided a more comprehensive theory of the interaction between two conducting plates which, among other things, handled the case of *finite* conductivity. This expression for the force has a Taylor series expansion in $(1/\sigma)$, where σ is the conductivity, and the leading order term — obtained in the limit of infinite conductivity — gave precisely Casimir's result. Again, this could be obtained as a result of direct electromagnetic coupling between the metals as long as a stochastic, fluctuating force is introduced. None of these, by itself, could be termed mysterious.

The mystery of the Casimir effect arises from the fact that it could be re-derived *without* really using the interaction between the constituents of the two metals. All the effect of these constituent charged particles is contained in a simple boundary condition on the metals in the infinite conductivity limit. The rest of the calculation uses the vacuum fluctuations of the electromagnetic field, its energy and a rather dubious subtraction scheme — however rigorously you deal with divergent series! — to get the same result. As I have explained, part of the mystery is in intuitively understanding why such different procedures lead to the same final expression.

The situation is complicated by another fact: While the force between two parallel metallic conductors is *attractive*, this need not be a general feature. For example, if you make two hemispherical bowls from conducting shells of matter and bring them close together, the Casimir force between them is actually repulsive; the geometry makes a lot of difference!

You can get Casimir effect without QFT mumbo-jumbo

But how come mumbo-jumbo gives the same result?

Plates attract but spheres repel!

Incredible, true ...

*... but still not
totally clear*

The whole phenomena is quite bewildering and if you are shaking your head in disbelief, I will not blame you! But the reality of this effect is beyond dispute and it has been derived from several different perspectives over years. The essential lesson is that the pattern of quantum fluctuations is sensitive to the boundary conditions we impose, both mathematically and practically. The ground state of the electromagnetic field in the presence of two parallel, conducting plates is quite different from the ground state in the absence of the plates. This alone is easy to understand because the ground state in the presence of the plates must ensure that, the field configurations which do not satisfy the boundary conditions at the plates, have zero probability for their existence. But what is rather curious is that this ground state has an energy which differs from that in the absence of the plates by a finite amount. There is no simple explanation for this fact, which makes the Casimir effect all the more fascinating.

Appendix: We will first derive an integral representation for the zeta function $\zeta(s)$ which allows analytic continuation for negative values of s . To do this, we begin with the integral representation for $\Gamma(s)$ and express it in the form:

$$\Gamma(s) = \int_0^\infty d\bar{t} \bar{t}^{s-1} e^{-\bar{t}} = n^s \int_0^\infty dt t^{s-1} e^{-nt}, \quad (19.37)$$

where we have set $\bar{t} = nt$ to arrive at the last expression. Summing the expression for $\Gamma(s)/n^s$ over all n , we get the result:

$$\Gamma(s) \sum_{n=1}^\infty \frac{1}{n^s} = \int_0^\infty dt \left(\frac{t^{s-1}}{e^t - 1} \right) = \Gamma(s) \zeta(s). \quad (19.38)$$

*The magic of
analytic
continuation*

Let us now consider the integral in the complex plane

$$I(s) \equiv \int_C dz \frac{z^{s-1}}{e^z - 1}, \quad (19.39)$$

over a contour consisting of the following paths: (i) Along the real line from $x = \infty$ to $x = \varepsilon$ where ε is an infinitesimal quantity; (ii) On a circle of radius ε around the origin going from $\theta = 0$ to $\theta = 2\pi$; (iii) Along the real line from ε to ∞ . Along the contour in (i) we get the contribution $-\Gamma(s)\zeta(s)$; it can be easily shown that the contribution along the circle will only pick up the residue at the origin. Finally, along (iii) one obtains the contribution $e^{2\pi is} \zeta(s) \Gamma(s)$. We therefore find that

$$\begin{aligned} \int_C dz \frac{z^{s-1}}{e^z - 1} &= \Gamma(s) \zeta(s) [e^{2\pi is} - 1] = \Gamma(s) \zeta(s) e^{i\pi s} (2i \sin(\pi s)) \\ &= \zeta(s) e^{i\pi s} (2\pi i) \frac{1}{\Gamma(1-s)}. \end{aligned} \quad (19.40)$$

In arriving at the last expression, we have used the standard identity

$$\Gamma(s)\Gamma(1-s) = \frac{\pi}{\sin \pi s} . \quad (19.41)$$

This allows us to express $\zeta(s)$ as a contour integral in the complex plane given by

$$\zeta(s) = e^{-i\pi s} \Gamma(1-s) \frac{1}{2\pi i} \int_C dz \frac{z^{s-1}}{e^z - 1} . \quad (19.42)$$

By studying the analytical properties of the right hand side, it is easy to show that this expression remains well defined for negative integral values of s . For example, when $s = -1$, the integrand in the contour integral can be expressed as $(1/z^3)[z/(e^z - 1)]$. The residue at the origin is therefore governed by the z^2 term in the power series expansion of $[z/(e^z - 1)]$, giving a contribution proportional to B_2 . Similarly, for $s = -3$, the integrand in the contour integral can be expressed as $(1/z^5)[z/(e^z - 1)]$. The residue at the origin is therefore governed by the z^4 term in the power series expansion of $[z/(e^z - 1)]$, giving a contribution proportional to B_4 . In fact, it is easy to show, putting all the factors together, that

Such tricks are used very often and are worth learning

$$\zeta(1-2k) = -\frac{B_{2k}}{2k} \quad \text{for } k = 1, 2, \dots . \quad (19.43)$$

This gives $\zeta(-1) = -B_2/2 = -1/12$, $\zeta(-3) = -B_4/4 = 1/120$ etc. These results provide an alternate way of giving meaning to the divergent series which occur in the computation of the Casimir effect.

Radiation: Caterpillar becomes Butterfly

20

The electric field of a charged particle at rest has two key properties: (i) It decreases as $(1/r^2)$ where r is the distance from the charged particle. (ii) It is directed radially outward from the position of the charged particle. Given the Coulomb field of a charge at rest, one can find the field of a charge moving with uniform velocity by a Lorentz transformation. This leads to the result that, if the charge moves with a *uniform* velocity \mathbf{v} , the field is given by

$$\mathbf{E} = \frac{q\mathbf{r}}{r^3} \frac{(1 - v^2/c^2)}{(1 - (v^2/c^2) \sin^2 \theta)^{3/2}}; \quad \mathbf{B} = \frac{1}{c} \mathbf{v} \times \mathbf{E}, \quad (20.1)$$

where θ is the angle between the direction of motion and the radius vector \mathbf{r} which has the components $(x - Vt, y, z)$. This, of course, looks more complicated but, as we would have expected, it still shares the two key properties with the fields produced by the static charge. The electric field falls as $(1/r^2)$ at large distances and it is radially directed from the *instantaneous position* of the charge.

Caterpillar: (1) $1/r^2$ fall-off; (2) points radially

The energy flux corresponding to the electromagnetic field, given by the Poynting vector, scales as the square of the electromagnetic field. If the field decreases as $(1/r^2)$, the energy flux will fall as $(1/r^4)$ and, since the area of a spherical surface scales as r^2 , the total energy flowing through a sphere at large distances from the charge falls as $r^2 \times (1/r^4) = (1/r^2)$. Therefore, one cannot transfer energy to large distances with this kind of field. This is understandable because such a transfer cannot take place in the rest frame of the charge — in which we only have a static Coulomb field — and since we expect physical processes to be Lorentz invariant, it should not happen for a charge moving with uniform velocity either.

But, when the charge is accelerating, something dramatic happens. The electric field picks up an additional term which falls only as $(1/r)$ at large distances. The change from the $(1/r^2)$ dependence to the $(1/r)$ dependence makes a huge difference! When the field falls as $(1/r)$ at large dis-

Butterfly: (1) $1/r$ fall-off; (2) points transversely

tances, the energy flux will fall as $(1/r^2)$ and the total energy flowing through a sphere at large distances from the charge is $r^2 \times (1/r^2)$ which is a constant! Therefore, the fields arising from an accelerated charge are capable of transmitting energy to large distances from the charge. Clearly, it would be nice to understand better how acceleration leads to such a shift from the $(1/r^2)$ to $(1/r)$ dependence — which changes the caterpillar to a butterfly.

There is also another peculiar feature that arises when the charge undergoes an accelerated motion. The Coulomb field of a charge at rest, and that of a charge moving with a uniform velocity, is radial. That is, the electric field vector in these cases points radially outward from the instantaneous position of the charge. But in the case of accelerated motion, the electric field picks up a transverse component which is *perpendicular* to the radial direction. Since a propagating electromagnetic plane wave, for example, will have an electric field that is transverse to the direction of propagation of the wave, this fact is crucial for identifying the field generated by the acceleration with the electromagnetic radiation.

*Let us forget the
textbook derivation!*

There is a remarkably elegant and simple way of understanding both these features [76]. This approach, originally due to J.J. Thomson [77], deserves to be more widely known and possibly could replace the rather unimaginative derivation using Lienard-Wiechert potentials in the classrooms! (Thomson's approach is discussed, for example, in [20] and also appears in the standard text books [78, 79]; unfortunately these textbooks create an impression that the result is valid only for non-relativistic motion.) I will describe this approach and its essential features.

*Position, velocity,
acceleration;
nothing else matters*

To begin with, let us recall a few elementary facts about Maxwell equations which connect the electromagnetic fields to the motion of the source. Since the electric field is $\mathbf{E} = -(1/c)(\partial\mathbf{A}/\partial t) - \nabla\phi$, we see that the electric field has a component which depends linearly on $(\partial\mathbf{A}/\partial t)$. It is also well known that the source for the vector potential \mathbf{A} is the current \mathbf{j} , in the sense that $\square\mathbf{A} \propto \mathbf{j}$. Therefore, $(\partial\mathbf{A}/\partial t)$ will have a source that depends on $(\partial\mathbf{j}/\partial t)$. Since \mathbf{j} is linear in the velocity of the charge, we conclude that the electric field will have a source term which is linear in the time derivative of the velocity, viz., the acceleration \mathbf{a} , but not on \dot{a} , \ddot{a} , ... etc.

An alternative way of understanding this result is as follows: A charge q moving with uniform velocity \mathbf{v} is equivalent to a current $\mathbf{j} = q\mathbf{v}$. This current will produce a magnetic field (in addition to the electric field) which scales in proportion to \mathbf{j} . If $\mathbf{a} = \dot{\mathbf{v}} \neq 0$, it will produce a non-zero $(\partial\mathbf{j}/\partial t)$ and hence a non-zero $(\partial\mathbf{B}/\partial t)$. Through Faraday's law, $(\partial\mathbf{B}/\partial t)$ will induce an electric field which scales as $(\partial\mathbf{j}/\partial t)$. (That is, if $(\partial\mathbf{j}/\partial t)$ changes by factor 2, the electric field will change by factor 2.) It follows that an accelerated charge will produce an electric field which is linear in $(\partial\mathbf{j}/\partial t) = q\mathbf{a}$. [This field, of course, is in addition to the usual Coulomb term which is independent of \mathbf{a} and falls as r^{-2} .]

Further, since the wave equation $\square \mathbf{A} \propto \mathbf{j}$ propagates information at the speed of light, we also know that the electric field at an event (t, \mathbf{x}) is determined entirely by the behaviour of the source at the event (t_R, \mathbf{x}') where $t - t_R = (1/c)|\mathbf{x} - \mathbf{x}'| \equiv (r/c)$. It is usual to call t_R as the “retarded time”.

Not now, but then

Before we do any sophisticated mathematics, let us try a bit of dimensional analysis to determine the electric field which arises from the acceleration. We know that the electric field has to be determined by the charge of the particle q , speed of light c , acceleration a and distance r (with a and r calculated at the retarded time.) In general, the field will also depend on the velocity of the particle at the retarded time but we will choose a Lorentz frame in which the charge was at rest *at the retarded time*, thereby eliminating any v dependence. We next use the fact that the electric field — which is linear in $\partial \mathbf{j} / \partial t$ — should be linear in both q and a , to write

Power of dimensional analysis

$$E = C(\theta) \frac{qa}{c^n r^m} = C(\theta) \left(\frac{q}{r^2} \right) \left(\frac{a}{c^n r^{m-2}} \right), \quad (20.2)$$

where C is a dimensionless factor, which can depend only on the angle θ between \mathbf{r} and \mathbf{a} , and n and m need to be determined. (Since $\mathbf{v} = 0$ in the instantaneous rest frame, the field cannot depend on the velocity.) From dimensional analysis, noting that E has the dimensions of q/r^2 it immediately follows that $(a/c^n r^{m-2})$ must be dimensionless, leading to $n = 2, m = 1$. So we get the result:

$$E = C(\theta) \frac{qa}{c^2 r}. \quad (20.3)$$

Thus, dimensional analysis plus the fact that \mathbf{E} must be linear in q and a , implies the r^{-1} dependence for the radiation term.

We get the a/r dependence

While this result shows why a term linear in acceleration will also have a $(1/r)$ dependence, it does not really tell us how exactly it comes about. Moreover, dimensional analysis cannot determine the nature of dimensionless function $C(\theta)$. The argument due to J.J. Thomson [77] does both of these in a rather neat way and I will now describe a slightly modified version of the same.

Let us consider a charged particle \mathcal{A} moving along some arbitrary trajectory $\mathbf{z}(t)$. We are interested in the electric field, say, produced at an event $\mathcal{P}(t, \mathbf{x})$ by this charge. Since the characteristics of the wave equation shows that information propagates at the speed of light from the source point to the field point, we already know that the field at \mathcal{P} will be determined by the properties of the trajectory at the retarded time t_R . Further, the electric field can only depend on the position $\mathbf{z}(t_R)$, velocity $\dot{\mathbf{z}}(t_R)$ and the acceleration $\ddot{\mathbf{z}}(t_R)$ at the retarded time but not on higher time derivatives. (This is because the source for electromagnetic field only involves up to the first time derivative of the current which is proportional to the acceleration.) We will now choose our Lorentz frame such that the

A series of clever choices

charge was at rest at the origin of the spacetime coordinates at the retarded time $t_R = 0$. Let the acceleration of the charge be $\mathbf{a} = \ddot{\mathbf{z}}(t_R)$ at this instant. We will rotate the coordinate system so that \mathbf{a} is along the x -axis.

We now consider *another* charged particle \mathcal{B} which was at rest, at the origin, from $t = -\infty$ to $t = 0$ and undergoes constant acceleration \mathbf{a} along the x -axis for a short time Δt . For $t > \Delta t$, it moves with constant velocity $v = a\Delta t$ along the x -axis. Let us study the electric field produced by this charge \mathcal{B} at some time $t \gg \Delta t$. Since Δt is arbitrarily small, we have $a\Delta t \ll c$ and we can use the non-relativistic approximation throughout. *Since the trajectory of \mathcal{B} matches identically in position, velocity and acceleration with the trajectory of \mathcal{A} , we are really interested in, it follows that both of them will produce identical electric fields at \mathcal{P} . This was the key insight of Thomson. As we shall see, the field produced by \mathcal{B} is fairly trivial to calculate and hence we can obtain the field due to \mathcal{A} .*

Introduce another charge you can control ...

... which produces the same field

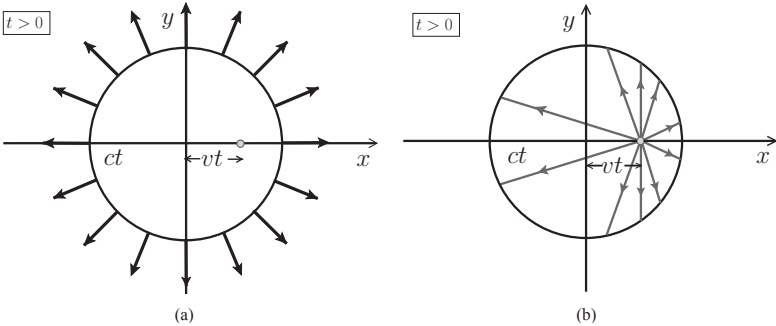


Fig. 20.1: A charge was at rest at the origin until $t = 0$ and was accelerated for a short amount of time δt after which it moves with a uniform velocity along the x -axis. The figure shows the electric field of the charged particle at some time $t > 0$. (a) The information that the charge was accelerated has not reached the region $r > ct$. In this region, the field is Coulombic and is radially outwards from the origin. (b) In the region $r < ct$, the field is that of a charge moving with uniform velocity. This field is radially outward from the instantaneous position of the charge.

The ‘news’, that the charge was accelerated at $t = 0$, could have only traveled up to a distance $r = ct$ in time t . Thus, at $r > ct$, the electric field should be that due to a charge located at the origin as shown in part (a) of **Fig. 20.1**:

$$\mathbf{E} = \frac{q}{r^2} \hat{\mathbf{r}} \quad (\text{for } r > ct) . \tag{20.4}$$

At $r \lesssim ct$, the field is that due to a charge moving with velocity v along the x -axis, given by Eq. (20.1). The key point is that this field is radially directed from the *instantaneous position* of the charge. When $v \ll c$, which is the situation we are interested in, this is again a Coulomb field radially directed from the *instantaneous* position of the charged particle (see part

Identify the two regions of Coulomb field

(b) of Fig. 20.1):

$$\mathbf{E} = \frac{q}{r'^2} \hat{\mathbf{r}}' \quad (\text{for } r < ct). \quad (20.5)$$

Around $r = ct$, there exists a small shell of thickness $(c\Delta t)$ in which neither result holds good. It is clear that the electric field in the transition region should interpolate between the two Coulomb fields. The crucial question is how we do this while ensuring that the flux of the electric field vector through any small box in this region vanishes, as it should in order to satisfy the Maxwell equations. As we shall see below, it turns out that this requires the field lines to appear somewhat like those shown in Fig. 20.2. We have concentrated on a single field line in Fig. 20.3 for clarity. One can explicitly work out this condition and prove that $\tan \theta = \gamma \tan \phi$ where $\gamma = (1 - v^2/c^2)^{-1/2}$. (It is done in detail in [78].) In the non-relativistic limit that we are considering, $\theta \approx \phi$ making the field lines parallel to each other in the inside and outside regions; that is, QP is parallel to RS. (This is easy to understand because the radial field is just the Coulomb field both in the outside and in the inside region. For the flux to be conserved, these two field lines should be parallel to each other.) What is really interesting is that we now need a piece of electric field line PR interpolating between the two Coulomb fields. This is clearly transverse to the radial direction and all that we need to do is to prove that its magnitude varies as $1/r$. Let us see how this comes about.

The shell of radiation

You should try this out!

Why is radiation field transverse?

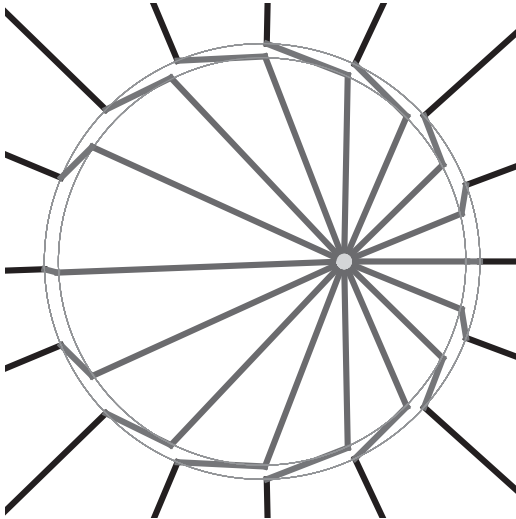


Fig. 20.2: Combining the field configurations at $r > ct$ and $r < ct$, shown in Fig. 20.1, requires the introduction of a transverse electric field in the transition region. Radiation arises from the necessity to connect two Coulomb fields in the regions $r > ct$ and $r < ct$ conserving the electric flux.

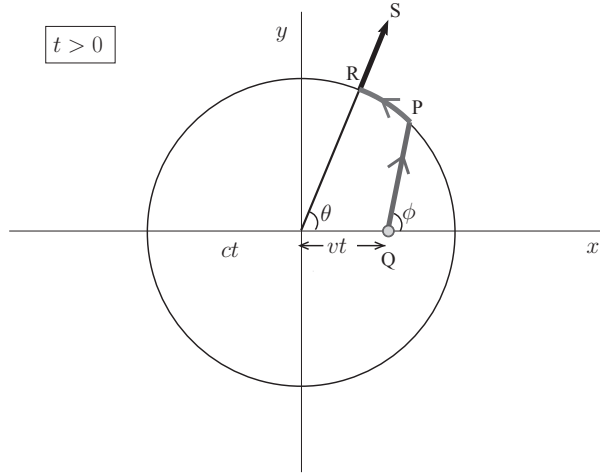


Fig. 20.3: One specific electric field line showing the way a transverse component is developed. The field line RS (in the region $r > ct$) is to be connected with the field line QP in the region $r < ct$ by the field line PR conserving the electric flux. This uniquely fixes the field line PR, which turns out to be the radiation field.

The situation is described in detail in Fig. 20.4 which is self-explanatory. Let E_{\parallel} and E_{\perp} be the magnitudes of the electric fields parallel and perpendicular to the direction \hat{r} . From the geometry, we have

$$\frac{E_{\perp}}{E_{\parallel}} = \frac{v_{\perp} t}{c \Delta t} . \quad (20.6)$$

But $v_{\perp} = a_{\perp} \Delta t$ and $t = (r/c)$, giving:

$$\frac{E_{\perp}}{E_{\parallel}} = \frac{(a_{\perp} \Delta t) (r/c)}{c \Delta t} = a_{\perp} \left(\frac{r}{c^2} \right) . \quad (20.7)$$

The value of E_{\parallel} can be determined by using Gauss' theorem to a small pill box, as shown in the small inset in Fig. 20.4. This gives $E_{\parallel} = E_r = (q/r^2)$; thus, we find that

$$E_{\perp} = a_{\perp} \left(\frac{r}{c^2} \right) \cdot \frac{q}{r^2} = \frac{q}{c^2} \left(\frac{a_{\perp}}{r} \right) . \quad (20.8)$$

This is the radiation field located in a shell at $r = ct$, which is propagating outward with a velocity c . The above argument clearly shows that the origin of the r^{-1} dependence lies in the necessity to interpolate between two Coulomb fields. We have thus determined the electric field generated due to the acceleration of the charge and have shown that it is transverse and falls as $(1/r)$!

Why does radiation field fall as $1/r$?

And why does it travel at the speed of light?

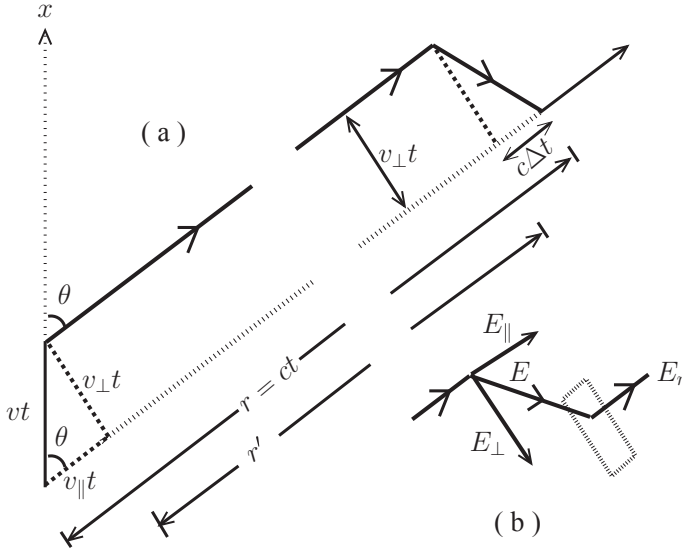


Fig. 20.4: (a) The electric field due to a charged particle which was accelerated for a small time interval Δt . For $t > \Delta t$, the particle is moving with a uniform non-relativistic velocity v along the x -axis. At $r > ct$, the field is that of a charge at rest in the origin. At $r < c(t - \Delta t)$, the field is directed towards the instantaneous position of the particle. The radiation field connects these two Coulomb fields in a small region of thickness $c\Delta t$. (b) Pill box construction to relate the normal component of the electric field around the radiation zone.

We can express this result more concisely in the vector notation as:

$$\mathbf{E}_{\text{rad}}(t, \mathbf{r}) = \frac{q}{c^2} \left[\frac{1}{r} \hat{\mathbf{n}} \times (\hat{\mathbf{n}} \times \mathbf{a}) \right]_{\text{ret}}, \quad (20.9)$$

where $\mathbf{n} = (\mathbf{r}/r)$ and the subscript “ret” implies that the expression in square brackets should be evaluated at $t' = t - r/c$. Comparison with Eq. (20.3) shows that $C(\theta) = \sin \theta$. The full electric field *in the frame in which the charge is instantaneously at rest*, is $\mathbf{E} = \mathbf{E}_{\text{coul}} + \mathbf{E}_{\text{rad}}$.

We emphasize that this result is *exact* in the Lorentz frame in which the charge was at rest at the retarded time. One does not have to make a non-relativistic “approximation” because $v = 0$ automatically takes care of it!. If we now make a Lorentz transformation to a frame in which the particle was moving with some velocity $\mathbf{v} = \dot{\mathbf{z}}(t_R)$ at the retarded time, then we can obtain the standard, fully relativistic expression with the velocity dependence. This is algebraically a little complicated because one needs to make a Lorentz transformation in an arbitrary direction since \mathbf{v} and \mathbf{a} will not — in general — be in the same direction.

Exact but in a special Lorentz frame

*The general result,
in full relativistic
glory*

Fortunately, there is an elegant way of doing this using 4-dimensional tensor notation which maintains manifest relativistic invariance. This will provide a complete relativistically invariant expression for the electromagnetic field of an arbitrarily moving charged particle without us ever having to mention the Lienard-Wiechert potential. Thus J.J.Thomson's idea is quite capable of giving us the complete solution to the problem. I outline this analysis; more details can be found in Ref. [80].

Consider a charge moving along an arbitrary trajectory $z^i(\tau)$ whose electromagnetic field $F^{ab}(x^i)$ at the observation point x^i is to be evaluated. We shall use units in which $c = 1$. The electromagnetic field tensor F^{ij} is then found by the relation $F^{ij} = \partial^i A^j - \partial^j A^i$ where A^i satisfies the equation $\square A^i = -4\pi J^i$ with J^i being the current. The $F^{0\alpha}$ terms give the components of the electric field, and the $F^{\alpha\beta}$ terms lead to the components of the magnetic field.

*Depends on
retarded time,
position, velocity
and acceleration*

We begin by noting that the electromagnetic field at the observation point x^i may depend only on the relative position $R^i = x^i - z^i(\tau)$, the velocity u^i , and the acceleration $a^i = du^i/d\tau$ of the charge, all evaluated at the retarded time τ_{ret} , but not on further derivatives of the trajectory. This result arises from the following: (a) Because electromagnetic signals propagate at the speed of light, the field at x^i is determined by the state of the source at an earlier position $z^i(\tau_{\text{ret}})$ which is related to x^i by a null line; that is, by the condition $R_i R^i = 0$. Of the two roots to this equation, we choose the retarded (causal) solution that satisfies the condition $R^0 > 0$. This condition determines the retarded time τ_{ret} . (b) Translational invariance implies that the field depends only on the relative position R^i of the charge with respect to the observation point (evaluated at the retarded time), and not on the absolute positions of the source or the observation point separately. (c) Because $\square A^i \sim J^i$, F^{ik} satisfies $\square F^{ik} \sim \partial^i J^k - \partial^k J^i$. Because J^i is at most linear in the velocity of the charge, $\partial^i J^k$ is at most linear in the acceleration, and no further derivatives of the trajectory can occur in the solution F^{ik} . Therefore, F^{ij} is a second rank antisymmetric tensor which is built from R^i , u^i , and a^i .

Some useful scalars

At this stage it is convenient to introduce the Lorentz invariant scalar $\ell = R_i u^i$ which, in the rest frame of the charge, reduces to:

$$\ell = R_i u^i = -R^0 = -|\mathbf{R}| \equiv -R, \quad (20.10)$$

where $(R^0)^2 = |\mathbf{R}|^2$ because of the condition $R_i R^i = 0$ and $R^0 > 0$ for the retarded solution. For simplicity, we will also define a four-vector n^i through the relation $R^i \equiv -\ell(n^i + u^i)$. It is easy to see that $n_k u^k = 0$, and $n_k n^k = 1$. The components of n^i are:

$$n^i = \left(-\frac{R}{\ell} - \gamma, -\frac{\mathbf{R}}{\ell} - \gamma \mathbf{v} \right), \quad (20.11)$$

which reduces, in the rest frame of the charge, to the unit spatial vector pointing from the charge to the field point: $n^i = (0, \mathbf{1})$. We will trade off the R^i dependence of F^{ij} for the n^i dependence and treat F^{ij} as a function of n^i , u^i , and a^i (instead of R^i , u^i , and a^i).

We next construct two four vectors E^i and B^i , defined as:

$$E^i = u_j F^{ij}, \quad B^i = \frac{1}{2} \epsilon^{ijkl} u_j F_{kl}, \quad (20.12)$$

*Electric and
Magnetic fields,
relativistic notation*

where ϵ^{ijkl} is the totally antisymmetric tensor in $D = 4$. The vectors E^i and B^i contain the same amount of information as F^{ij} as can be seen by the explicit expression for the latter in terms of the former:

$$F^{ij} = u^i E^j - E^i u^j - \epsilon^{ij}_{\quad kl} u^k B^l, \quad (20.13)$$

which can be easily verified by direct substitution of Eq. (20.13) into Eq. (20.12) and the use of the identities $u_j E^j = 0$ and $u_j B^j = 0$. These identities also show that E^i and B^i are both orthogonal to u^i , and hence, in a given reference frame, they contain only three independent components as required.

The four vectors E^i and B^i have direct physical interpretations and represent the electric and magnetic fields in the instantaneous rest frame of the charge with four velocity u^i . In this frame, $u^j = (1, \mathbf{0})$ so that only the component u_0 contributes, and $E^i = u_0 F^{i0} = (0, F^{0\alpha})$, because F^{ij} is antisymmetric. Hence, the spatial components of $E^i = u_j F^{ij}$ correctly represent the components of the electric field in the instantaneous rest frame of the charge. Similarly, in the instantaneous rest frame, only the component u_0 contributes to B^i . Because ϵ^{ijkl} is completely antisymmetric, the time component of B^i vanishes in this frame. We see that the spatial components of B^i are given by $F^{\alpha\beta}$ where $\alpha, \beta = 1, 2, \text{ or } 3$. Hence the spatial components of B^i lead to the correct values of the magnetic field components in the rest frame.

However, we already know the form of the electromagnetic field in the instantaneous rest frame from Thomson's argument: The electric field is given by Eq. (20.9) with the magnetic field given by $\mathbf{B} = \hat{\mathbf{n}} \times \mathbf{E}$. In the rest frame, we have $n^i = (0, \mathbf{1})$, $u^j = (1, \mathbf{0})$, $a_i = (0, \mathbf{a})$ and $R_i u^i = \ell$. Using these, it is easy to see that Thomson's electromagnetic fields can be expressed in four-dimensional notation as:

*The trick: Write
Thomson's result in
relativistic nota-
tion ...*

$$\mathcal{E}^i = \frac{q}{\ell^2} n^i + \frac{q}{\ell} [a^i - n^i (n_k a^k)]; \quad \mathcal{B}^i = \frac{q}{\ell} \epsilon^{ijkl} u_j n_k a_l. \quad (20.14)$$

In fact, this completely solves the problem and provides the electric and magnetic fields in any Lorentz frame. But if we are interested in determining F^{ij} (since this is the usual quantity used in relativistic physics), we can do that by substituting the four-dimensional generalized fields de-

*... if you want the
textbook result*

rived from the Thomson expression, namely \mathcal{E}^i and \mathcal{B}^i , for E^i and B^i respectively in Eq. (20.13) and obtain the explicit expression for F^{ij} . If we substitute Eq. (20.14) into Eq. (20.13), we obtain

$$F^{ij} = \frac{q}{\ell^2} u^{[i} n^{j]} - \frac{q}{\ell} a^{[i} u^{j]} + \frac{q}{\ell} (n_k a^k) n^{[i} u^{j]} - \frac{q}{\ell} \varepsilon^{ij}_{kl} \varepsilon^{lpqr} u^k n_q a_r u_p . \quad (20.15)$$

To evaluate the expression $\varepsilon^{ij}_{kl} \varepsilon^{lpqr} u^k n_q a_r u_p$ we use the identity

$$\varepsilon_{ijkl} \varepsilon^{lpqr} = -[\delta_i^p (\delta_j^r \delta_k^q - \delta_j^q \delta_k^r) - \delta_i^r (\delta_j^p \delta_k^q - \delta_j^q \delta_k^p) + \delta_i^q (\delta_j^p \delta_k^r - \delta_j^r \delta_k^p)] . \quad (20.16)$$

We lower the indices i and j in Eq. (20.15) and then use Eq. (20.16) to obtain the expression for F_{ij} :

$$F_{ij} = \frac{q}{\ell^2} u_{[i} n_{j]} - \frac{q}{\ell} a_{[i} u_{j]} + \frac{q}{\ell} (n_k a^k) n_{[i} u_{j]} + \frac{q}{\ell} n_{[i} a_{j]} . \quad (20.17)$$

The final result is known in the literature. It is obtained by integrating the Maxwell equations in a four-dimensional notation, and differentiating the resultant Lienard-Wiechert potentials A^j with respect to x^i . The present approach is significantly more elegant and simpler; if you do not believe me, try differentiating the Lienard-Wiechert potential!

The importance of the $(1/r)$ field, of course, is that it allows propagation of energy in the form of radiation to large distances. The amount of energy radiated by the system per unit time is given by the Larmor formula

$$\frac{dE}{dt} = \frac{2}{3} \frac{q^2}{c^3} a^2 , \quad (20.18)$$

which can be easily obtained from the results obtained above. We will conclude this chapter by discussing some interesting features related to this formula.

The point of historical importance is related to the radiation emitted by charges in circular motion. Obviously, a single charge, going around a circle of radius r with speed v , will have an acceleration $a = v^2/r$ and the radiated power will vary as $(v/c)^4$. But suppose we have two charges located at diametrically opposite points in a circle with both moving at the same speed around the circle. In this case, it is easy to show from symmetry that the dipole moment vanishes and the radiation has to come from the variation of the quadrupole moment; it will now be proportional to $(v/c)^6$. Similarly, if we think of three charged particles (all having the same charge) located 120 degrees apart in a circle, undergoing uniform circular motion, then we get only octupole radiation. In general, if we have N charged particles evenly spaced in a ring all co-moving in a circular orbit, then the radiation will be down by a factor $(v/c)^{2(n+1)}$. In fact, this is the reason we would often ignore any radiative field from a steady current going around in a circular orbit.

This problem was posed and the result was first obtained, again, by J.J. Thomson [81]. He was trying to explain the fact that electrons in the atoms do not radiate and he used this calculations to support his model that the electronic charge in an atom must be smoothly distributed. This was followed up by Schott [82] by an extensive analysis re-deriving the results. Curiously enough, all these were soon forgotten and were, in a way, re-invented around late 1940s when one wanted to study relativistic electrons in particle accelerators (see, for e.g., Ref. [83, 84]).

The reason Thomson worried about all these

Box 20.1: Radiation and Gauss law

Several textbooks introduce the Gauss law $\nabla \cdot \mathbf{E} = 4\pi\rho$ fairly early on in electrostatics and connect it up with Coulomb's law. The integral form of the Gauss law:

$$\int \mathbf{E}(t, \mathbf{x}) \cdot \mathbf{n} dA = 4\pi Q(t), \quad (20.19)$$

when applied to a point charge at rest immediately tells you that the electric field falls as $(1/r^2)$. Since the surface area of a sphere increases as r^2 , the above relation immediately follows.

After having done a fair amount of electrostatics, the text books will describe radiation fields in a later chapter and never revisit the Gauss law. But if the Gauss law is tied to $(1/r^2)$ electric field and the radiation field has a $(1/r)$ component, does it mean that we can't use Gauss law in general? The issue is further complicated by the fact that the radiation fields depend on retarded time, while the Gauss law relates the electric field at time t to the charge distribution at the same time t . *No retardation!* Since many students seem to be associating the Gauss law with the Coulomb law and $(1/r^2)$, it is worth clarifying this point.

The Gauss law, being one of Maxwell's equation, is universally valid and is definitely valid for the radiation field as well. Figure 20.5 illustrates this dramatically. In a region of space, at a given time t , there are a set of charged particles ($q_1, q_2, q_3 \dots$ are shown explicitly) in arbitrary, accelerated states of motion. Their trajectories are indicated in the picture and the black dots indicate their positions at a given instant of time $t = t_0$. The fields produced by the charged particle are quite different from the Coulomb $(1/r^2)$ field and include the radiative component.

Gauss' law is applicable to radiation field as well

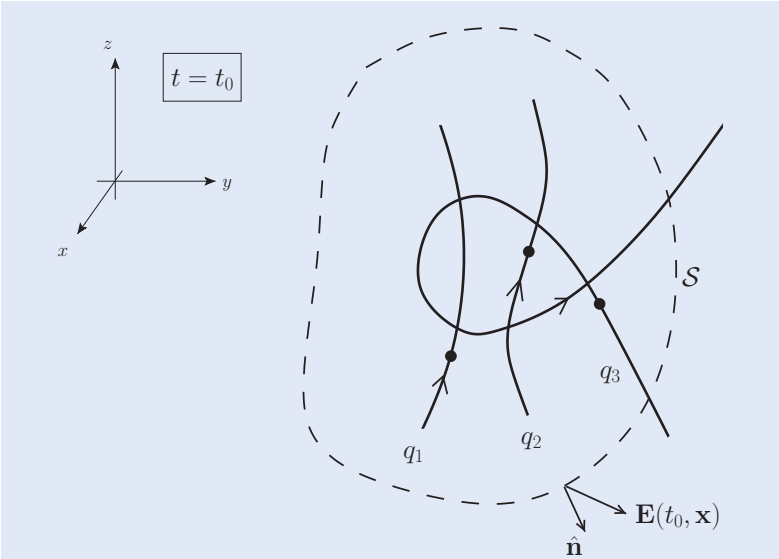


Fig. 20.5: Three charged particles are moving in space along the trajectories shown in the figure. At some time $t = t_0$, all the three charges happen to be inside a compact region of 3-space enclosed by the surface \mathcal{S} . Their positions are indicated by black dots. The electric field on \mathcal{S} is determined by the position, velocity and acceleration of these charged particles at respective retarded times when they might not have been inside \mathcal{S} . The field produced by the charges will also involve both Coulomb and radiation components. Nevertheless, the flux of the electric field through \mathcal{S} at time $t = t_0$ is precisely equal to, the total charge contained inside \mathcal{S} at the instant $t = t_0$ because of the Gauss law. Thus, Gauss law incorporates the radiation field and the retardation effect in a subtle manner.

Make sure you understand this

Let \mathcal{S} be a two dimensional compact surface as indicated in the figure by a broken line. The flux of the electric field $\mathbf{E}(t_0, \mathbf{x})$ at time $t = t_0$ through the surface \mathcal{S} will be precisely equal to $4\pi(q_1 + q_2 + q_3)$ for the situation shown in the figure. The three charges are inside \mathcal{S} at $t = t_0$ but they could have been outside at the respective retarded times. If there are other charges *outside* \mathcal{S} , they will *all* contribute to $\mathbf{E}(t_0, \mathbf{x})$ but not to the total charge count. And, as I have emphasized several times, the fields need not be purely Coulombic.

The magic of it

When you think about it, Gauss' law is quite fascinating. In fact, if one postulates a generally covariant version of the Gauss law to be valid for all observers in all states of motion, one can obtain all the Maxwell equations from it. It does not advertise special relativity, retardation effects, wave propagation and all that stuff but quietly recognizes them when considered as a part of the Maxwell equations!

You know that a blackbody cavity, kept at a given temperature T , will be filled with electromagnetic radiation of a particular spectral form, viz., the Planck spectrum. This is one case where we could have thought of the radiation either as fluctuating electric and magnetic fields, or as a bunch of photons. How does this dual role manifest itself? For example, if a charged particle interacts with the blackbody radiation, do we get the same results when we treat the radiation as fluctuating electromagnetic fields or as photons? We will try to understand this equivalence in some simple contexts in this chapter [20].

Blackbody radiation: Photons or Field?

To begin with, it is interesting to note that the blackbody radiation, *by itself*, exhibits both particle and wave nature. To see this, let us compute the energy fluctuations of the blackbody radiation. For a system in thermodynamic equilibrium at temperature T , with $\beta \equiv (kT)^{-1}$, the mean energy \bar{E} is given by

$$\bar{E} = \frac{\sum E e^{-\beta E}}{\sum e^{-\beta E}} = Z^{-1} \sum E e^{-\beta E} = -\frac{1}{Z} \frac{\partial Z}{\partial \beta}; \quad Z \equiv \sum e^{-\beta E}, \quad (21.1)$$

where Z is the partition function. Differentiating once again, we get an expression for mean square fluctuation in energy:

$$-\frac{\partial \bar{E}}{\partial \beta} = \frac{1}{Z} \frac{\partial^2 Z}{\partial \beta^2} - \left(\frac{1}{Z} \frac{\partial Z}{\partial \beta} \right)^2 = \langle E^2 \rangle - \bar{E}^2 = (\Delta E)^2. \quad (21.2)$$

In the case of blackbody radiation with $\bar{E} = \hbar \omega (e^{\beta \hbar \omega} - 1)^{-1}$, direct differentiation gives

$$(\Delta n)^2 \equiv \left(\frac{\Delta E}{\hbar \omega} \right)^2 = \left(\frac{\bar{E}}{\hbar \omega} \right)^2 + \left(\frac{\bar{E}}{\hbar \omega} \right) = \bar{n}^2 + \bar{n}, \quad (21.3)$$

where $\bar{n} \equiv (\bar{E}/\hbar \omega)$ is the mean number of photons with frequency ω and Δn is the fluctuation in this number. *It is both!*

Curiously enough, these two terms in Eq. (21.3) represent the fluctuations which will arise when we think of the system as made of waves or particles. If photons were to be interpreted as particles, then one would expect $(\Delta n)^2 \simeq \bar{n}$, giving the usual Poisson fluctuations $(\Delta n/n) \simeq n^{-1/2}$. For this to occur we will need $\bar{n} \gg \bar{n}^2$; that is, $\bar{n} \ll 1$, which happens for $\beta \hbar \omega \gg 1$. On the other hand, if $\beta \hbar \omega \ll 1$, we have $\bar{n} \gg 1$ and we get $(\Delta n)^2 \simeq \bar{n}^2$ which characterizes the wave-like fluctuations. In these two limits, given by $\hbar \omega \ll kT$ and $\hbar \omega \gg kT$, the expression for $n(\omega)$ itself has simple asymptotic behaviour consistent with the above interpretation.

*All as expected
and nice*

When $\hbar \omega \ll kT$, we are in the long wavelength, classical regime of the radiation. Equipartition of energy suggests that each mode (having two polarization states) should have energy $\epsilon_\omega = 2 \times (kT/2) = (kT)$ or $n_\omega = (\epsilon_\omega/\hbar \omega) = (kT/\hbar \omega)$. This is what we get from the Planck spectrum.

When $\hbar \omega \gg kT$, we are in the regime in which photons behave as particles. In that case, we expect $n_\omega = \exp(-\hbar \omega/kT)$ based on Boltzmann statistics. (In this limit, $n_\omega \ll 1$ and quantum statistical effects are ignorable; hence we get Boltzmann statistics rather than Bose-Einstein statistics.) Again, this is what we obtain from Planck spectrum. Thus, one may think of blackbody radiation as made of photons when $\hbar \omega \gg kT$ and as made of waves when $\hbar \omega \ll kT$.

After this warm up, let us consider a more complicated situation when the radiation field is not isolated but interacts with charged particles.

Consider a gas of electrons at temperature T_e , interacting with a distribution of photons with mean energy $\langle E \rangle$. Assume that $\langle E \rangle \ll mc^2$ and $kT_e \ll mc^2$. During the scattering, energy is exchanged between electrons and photons. We are interested in computing the net energy transfer between the charged particle and the photons. Obviously, we expect the net energy transfer $\langle \Delta E \rangle$ from the photons to the electrons to be positive if the average energy $\langle E \rangle$ of the photons is much larger than the thermal energy of the electrons; on the other hand, if $\langle E \rangle \ll k_B T_e$ we expect $\langle \Delta E \rangle$ to be negative with the photons getting the energy from the electrons. This is an interesting situation (which happens to be of considerable astrophysical importance) which we will analyse from different angles.

*The sheer power
of Taylor series
expansion!*

I will first show how the result can be obtained by a rather cute trick. Let the energy transferred from the photons to the electrons be ΔE on the average. (We shall omit the symbol $\langle \rangle$ for simplicity of notation.) Since $E \ll mc^2$ and $kT_e \ll mc^2$, we can expand $(\Delta E/mc^2)$ in a double Taylor series in (E/mc^2) and (kT_e/mc^2) , retaining up to quadratic order:

$$\begin{aligned} \frac{\Delta E}{mc^2} = & c_1 + c_2 \left(\frac{E}{mc^2} \right) + c_3 \left(\frac{kT_e}{mc^2} \right) + c_4 \left(\frac{E}{mc^2} \right)^2 + c_5 \left(\frac{E}{mc^2} \right) \left(\frac{kT_e}{mc^2} \right) \\ & + c_6 \left(\frac{kT_e}{mc^2} \right)^2 + \dots \end{aligned} \quad (21.4)$$

The coefficients (c_1, \dots, c_6) can be fixed by the following arguments:

(i) Since $\Delta E = 0$ for $T_e = E = 0$, we must have $c_1 = 0$.

(ii) Consider next the scattering of a photon with electrons at rest. This will correspond to $T_e = 0$ and $E \neq 0$. If the scattering angle is θ , the standard result of Compton scattering tells you that the wavelength of the photon changes by *Just textbook*

$$\Delta\lambda = \left(\frac{h}{mc}\right)(1 - \cos\theta). \quad (21.5)$$

Such scattering of a photon, by an electron at rest, is symmetric in the forward — backward directions and the mean fractional change in the frequency is $(\Delta\omega/\omega) = -(\Delta\lambda/\lambda) = -(\hbar\omega/mc^2)$; so the average energy transfer to the electrons is $\Delta E = E^2/mc^2$. This implies that $c_2 = 0$ and $c_4 = 1$.

(iii) If $E = 0$ and $T_e \neq 0$ the photon has zero energy and nothing should happen; hence $c_3 = c_6 = 0$. So, our expression reduces to

$$\frac{\Delta E}{mc^2} = \left(\frac{E}{mc^2}\right)^2 + c_5 \left(\frac{E}{mc^2}\right) \left(\frac{kT_e}{mc^2}\right)^2. \quad (21.6)$$

(iv) To fix c_5 , which is the really non-trivial coefficient, we can consider the following thought experiment. Suppose there is a very *dilute* gas of photons at the *same* temperature as the electrons. Then the number density $n(E)$ of photons is given by the Boltzmann limit of the Planck distribution: *For this you need a trick*

$$n(E)dE \propto E^2 \left(e^{\beta E} - 1\right)^{-1} dE \propto E^2 \exp(-E/kT_e) dE. \quad (21.7)$$

In this case, since the temperatures are the same, we expect the net energy transfer between the electrons and the photons to vanish. That is, we demand

$$0 = \int_0^\infty dE n(E) \Delta E, \quad (21.8)$$

in this situation. Substituting for ΔE and $n(E)$ from Eq. (21.6) and Eq. (21.7), one can easily show that $(4kT_e + c_5 kT_e) = 0$ or $c_5 = -4$. Hence, we get the final result

$$\frac{\Delta E}{E} = \frac{(E - 4kT_e)}{mc^2}. \quad (21.9)$$

One may say that, in a typical collision between an electron and photon, the electron energy changes by (E^2/mc^2) and the photon energy changes by $(4kT_e/mc^2)E$.

Let us explore this equation a bit more closely. From Compton scattering, we know that the average energy lost by the photon per collision is

given by

$$\langle \Delta \varepsilon \rangle = - \left(\frac{\hbar \omega_i}{m_e c^2} \right) \hbar \omega_i . \quad (21.10)$$

How does the radiation gain energy

Comparing with the result obtained above, we conclude that the mean fractional energy gained by the photon in one collision must be about $4k_B T_e / m_e c^2$. How do we interpret the generation of these *additional* photons, which, — classically — corresponds to radiation of electromagnetic waves? Why are the charges radiating? This turns out to be a bit more non-trivial and is related to a radiation drag force felt by the charged particle in a photon gas. So we first need to obtain the expressions for these. We will approach the problem step-by-step.

Radiated four-momentum

We begin by finding the relativistic analog of the Larmor formula for the radiation, which is given by (see Eq. (20.18)):

$$d\mathcal{E} = \frac{2}{3} \frac{q^2}{c^3} a^2(t') dt , \quad (21.11)$$

where t' is the retarded time. Let us choose an instantaneous rest frame for the charge in which this non-relativistic formula is valid at $t = t'$. Because of symmetry, the net momentum radiated, $d\mathbf{P}$, will vanish in this instantaneous rest frame. Clearly this result should be valid even for relativistic motion, if we can rewrite it in an invariant manner. If a^i is the four-acceleration, then $a^2/c^4 = a^i a_i$ in the instantaneous rest frame of the charge. So we can express Eq. (21.11), as well as the condition $d\mathbf{P} = 0$, in the form

$$dP^k = \frac{2}{3} \frac{q^2}{c} (a^i a_i) dx^k = \frac{2}{3} \frac{q^2}{c} (a^i a_i) u^k ds , \quad (21.12)$$

where dP^k is the four-momentum radiated by the particle during the proper time interval ds . Being relativistically invariant, this result is true for arbitrary velocities.

This radiation leads to a damping force on the particle which, in the fully relativistic case is given by a four force g^i (see Appendix for derivation):

$$g^i = \left(\frac{2q^2}{3} \right) \left[\frac{d^2 u^i}{ds^2} + u^i u^k \frac{d^2 u_k}{ds^2} \right] = \frac{2}{3} q^2 \left[\frac{d^2 u^i}{ds^2} - u^i (a^k a_k) \right] . \quad (21.13)$$

Special, but useful, case

These expressions are valid irrespective of the nature of the source which is accelerating the particle. But, in most contexts, this acceleration will be produced by an externally specified electromagnetic field. If this electromagnetic field is represented by the field tensor F^{ik} (which we assume to be a constant for the sake of simplicity), then we have:

$$a^i = \left(\frac{q}{m} \right) F^i_k u^k ; \quad \frac{da^i}{ds} = \left(\frac{q}{m} \right)^2 F^i_k F^k_j u^j . \quad (21.14)$$

Substituting these expressions in Eq. (21.13), and rearranging the terms, we get

$$g^i = -\frac{2}{3} \left(\frac{q^2}{m} \right)^2 \left[\left(F^{ka} F_{kj} \right) u_a u^j u^i + F^{ki} F_{kj} u^j \right]. \quad (21.15)$$

This expression can be written in a much nicer form in terms of the energy-momentum tensor for the electromagnetic field:

$$4\pi T_{bc} = F_{ab} F^a{}_c - \frac{1}{4} F_{mn} F^{mn} g_{bc}. \quad (21.16)$$

Using this expression, we can write

$$F^{il} F_{kl} = F^{li} F_{lk} = (4\pi) T^i_k + \frac{1}{4} \delta^i_k \left(F_{ab} F^{ab} \right). \quad (21.17)$$

Now we can express g^i in terms of T^{ab} alone without F_{ik} appearing explicitly. Note that, when we use Eq. (21.17) in Eq. (21.15) the term involving $F^2 = F_{ab} F^{ab}$ cancels out. Therefore

This is neat

$$\begin{aligned} g^i &= \frac{8\pi}{3} \left(\frac{q^2}{m} \right)^2 \left[T^{ij} u_j - \left(T^{ab} u_a u_b \right) u^i \right] \\ &= \left(\frac{\sigma_T}{c} \right) \left[T^{ij} u_j - \left(T^{ab} u_a u_b \right) u^i \right], \end{aligned} \quad (21.18)$$

where $\sigma_T = (8\pi/3) (q^2/mc^2)^2$ is the Thomson cross-section. This is a nice relation which expresses the radiation reaction force in terms of the energy-momentum tensor of the electromagnetic field which is accelerating the charged particle.

As a simple application of this result, consider the humble phenomenon of Thomson scattering. When an electromagnetic wave hits a charged particle, it makes the particle oscillate and radiate. The radiation will exert a damping force on the particle. In a frame in which the charge is at rest, $u^i = (1, 0, 0, 0)$ and $g^i = (\gamma \mathbf{f} \cdot \mathbf{v}, \gamma \mathbf{f}) = (0, \mathbf{f})$. From Eq. (21.18), we get:

Example: Thomson scattering

$$g^i = \sigma_T \left[T^{i0} - T^{00} u^i \right] = (0, \sigma_T u \hat{\mathbf{n}}), \quad (21.19)$$

which is a standard result.

For a more non-trivial example, let us come back to the problem of charged particles interacting with a radiation field with energy density U_{rad} . Using $T^{ab} = U_{\text{rad}} \text{dia} (1, 1/3, 1/3, 1/3)$ for an isotropic radiation bath and $u^i = (\gamma, \gamma \mathbf{v})$ we get

Back to the original problem

$$T^{ab} u_a u_b = U_{\text{rad}} \gamma^2 \left(1 + \frac{1}{3} v^2 \right); \quad T^{ab} u_b = \left(U_{\text{rad}} \gamma, -\frac{1}{3} U_{\text{rad}} \gamma \mathbf{v} \right). \quad (21.20)$$

This gives, on using Eq. (21.18),

$$g^i = \left(-\frac{4}{3} \sigma_T U_{\text{rad}} \gamma^3 v^2, -\frac{4}{3} \sigma_T U_{\text{rad}} \gamma^3 \mathbf{v} \right) = (\gamma \mathbf{f} \cdot \mathbf{v}, \gamma \mathbf{f}). \quad (21.21)$$

Comparing, we get

$$\mathbf{f} = -\frac{4}{3} \sigma_T U_{\text{rad}} \gamma^2 \left(\frac{\mathbf{v}}{c} \right); \quad -\mathbf{f} \cdot \mathbf{v} = \frac{4}{3} \sigma_T U_{\text{rad}} \gamma^2 \left(\frac{v^2}{c^2} \right) c, \quad (21.22)$$

where we have re-introduced the c -factor. This result is valid for *any* radiation field with energy density U_{rad} . The work done by this drag force is given by the second relation in Eq. (21.22).

But this should be equal to the net power radiated by the electron! In other words, this is the addition of energy to the photon field due to the energy radiated by the electrons. The mean number of photons scattered per second is $N_c = (\sigma_T c n_{\text{rad}}) = (\sigma_T c U_{\text{rad}} / \hbar \omega_i)$ where $\hbar \omega_i$ is the average energy of the photon defined by $\hbar \omega_i = (U_{\text{rad}} / n_{\text{rad}})$. Hence the average energy gained by the photon in one collision is

$$\langle \Delta E \rangle = \frac{4}{3} \gamma^2 \left(\frac{v}{c} \right)^2 \hbar \omega_i = \frac{4}{3} \gamma^2 \left(\frac{v}{c} \right)^2 \langle E \rangle. \quad (21.23)$$

In the relativistic limit, $\langle \Delta E / E \rangle \simeq (4/3) \gamma^2 \gg 1$, and this process can be a source of high energy photons. When $v \ll c$, the energy gain by photons per collision is $\langle \Delta E / E \rangle \simeq (4kT_e / m_e c^2)$. This is precisely the result obtained earlier in Eq. (21.10). So when we think of charged particles interacting with radiation field made of photons, everything works out fine.

But the thermal bath can also be thought of as made up of fluctuating electromagnetic fields with no mention of photons. How can we account for the increase in the energy of the thermal bath when it interacts with the charged particles? Let us see how this comes about.

When a charged particle and a photon scatter off each other with the photon gaining energy, we do not bat an eyelid; we think of this process to be somewhat akin to two billiard balls colliding with each other with one gaining the energy lost by the other. But the addition of energy to a large bunch of photons is equivalent to the increase in the radiation field when we look at it in the wave picture. Such radiation can only come from the acceleration of charged particles. What is the source of acceleration of a charged particle kept inside a blackbody cavity? It has to be the fluctuating electromagnetic field when we view everything in the wave perspective. The fluctuating acceleration of a charged particle in the random electromagnetic field of the blackbody cavity has to produce precisely the correct amount of radiation emission as one would have obtained by thinking everything through in terms of photons.

*Everything is fine
in photon picture*

The real issue

A thermal bath of photons is equivalent to a random superposition of electromagnetic radiation with $\langle E^2/4\pi \rangle = \langle B^2/4\pi \rangle = aT^4$ at any location. If the charge is not moving, then there is no net flux hitting the charge and there is no drag force. Suppose the charge is moving with velocity \mathbf{v} , in a frame S in which radiation is isotropic. We will now transform to a frame S' in which the charge is at rest. The energy flux in S' along the x -axis is

*Introduce
fluctuating
EM field*

$$\begin{aligned} T'^{0x} &= \gamma^2 \left[\left(1 + \frac{v^2}{c^2} \right) T^{0x} - \frac{v_x}{c} (T^{00} + T^{xx}) \right] \\ &= -\frac{v_x}{c} \gamma^2 (aT^4) \left(1 + \frac{1}{3} \right) = -\frac{4}{3} (aT^4) \left(\frac{v_x}{c} \right) \gamma^2. \end{aligned} \quad (21.24)$$

We have used the facts $T^{0x} = 0, T^{00} = aT^4, T^{xx} = (1/3)aT^4$. From Eq. (21.19), we find that

$$\mathbf{f}_{\text{drag}} = -(4/3)U_{\text{rad}}\gamma^2 (\mathbf{v}/c) \cong -\frac{4}{3}\sigma_T (aT^4) \left(\frac{\mathbf{v}}{c} \right), \quad (21.25)$$

which is precisely the result we obtained earlier!

We can also obtain the power radiated by the charged particles directly using the notion of just electromagnetic fields. To do this, we will first re-express the Larmor formula for the power radiated in terms of the electric and magnetic fields which produce the acceleration. Consider a frame S in which the particle has a velocity \mathbf{v} and acceleration \mathbf{a} . We now make a Lorentz transformation to a frame S' in which the charge is instantaneously at rest. In this frame:

*An alternative
derivation*

$$\mathbf{E}'_{\parallel} = \mathbf{E}_{\parallel}, \quad \mathbf{E}'_{\perp} = \gamma(\mathbf{E}_{\perp} + \mathbf{v} \times \mathbf{B}), \quad (21.26)$$

and the acceleration is $\mathbf{a}' = (q/m)\mathbf{E}'$. (We have set $c = 1$ to simplify the expressions.) Hence the instantaneous power radiated is

$$\begin{aligned} \frac{2}{3}q^2a'^2 &= \frac{2}{3}\frac{q^4}{m^2} \left[\mathbf{E}_{\parallel}^2 + \gamma^2 (\mathbf{E}_{\perp} + \mathbf{v} \times \mathbf{B})^2 \right] \\ &= \frac{2}{3}\frac{q^4}{m^2} \left[\mathbf{E}_{\parallel}^2 + \gamma^2 (\mathbf{E} + \mathbf{v} \times \mathbf{B} - \mathbf{E}_{\parallel})^2 \right] \\ &= \frac{2}{3}\frac{q^4}{m^2} \left[\gamma^2 (\mathbf{E} + \mathbf{v} \times \mathbf{B})^2 - \gamma^2 E_{\parallel}^2 v^2 \right]. \end{aligned} \quad (21.27)$$

In arriving at the last equation we have used the relations $\mathbf{E} \cdot \mathbf{E}_{\parallel} = E_{\parallel}^2$ and $\mathbf{E}_{\parallel} \cdot (\mathbf{v} \times \mathbf{B}) = 0$. Writing $E_{\parallel}^2 v^2 = (\mathbf{E} \cdot \mathbf{v})^2$, we get

$$\Delta \mathcal{E} = \frac{2}{3} \left(\frac{q^4}{m^2} \right) \gamma^2 \left[(\mathbf{E} + \mathbf{v} \times \mathbf{B})^2 - (\mathbf{E} \cdot \mathbf{v})^2 \right] \Delta t. \quad (21.28)$$

We next treat the radiation field as equivalent to an electromagnetic field with $\langle (E^2/8\pi) \rangle = \langle (B^2/8\pi) \rangle = (U_{\text{rad}}/2)$ with \mathbf{E} and \mathbf{B} randomly fluctuating around zero mean. In this case, we can again use Eq. (21.28) and average over \mathbf{E} and \mathbf{B} to obtain the net power. Now,

$$\mathcal{Q} \equiv \langle (\mathbf{E} + \mathbf{v} \times \mathbf{B})^2 - (\mathbf{E} \cdot \mathbf{v})^2 \rangle = \langle E^2 - (\mathbf{E} \cdot \mathbf{v})^2 \rangle + \langle (\mathbf{v} \times \mathbf{B})^2 \rangle, \quad (21.29)$$

since $\langle \mathbf{E} \cdot (\mathbf{v} \times \mathbf{B}) \rangle = \langle \mathbf{v} \cdot (\mathbf{B} \times \mathbf{E}) \rangle = 0$ due to random orientation of \mathbf{v} with respect to $(\mathbf{E} \times \mathbf{B})$. Using the relation

$$\langle E^2 - (\mathbf{E} \cdot \mathbf{v})^2 \rangle = \langle E^2 \rangle - E^2 v^2 \langle \cos^2 \theta \rangle = E^2 (1 - v^2/3) \quad (21.30)$$

and

$$\langle (\mathbf{v} \times \mathbf{B})^2 \rangle = \langle \mathbf{v} \cdot [\mathbf{v} B^2 - \mathbf{B}(\mathbf{v} \cdot \mathbf{B})] \rangle = v^2 B^2 - v^2 B^2/3 = (2/3) \beta^2 B^2, \quad (21.31)$$

we get $\mathcal{Q} = E^2 (1 - v^2/3) + (2/3) v^2 B^2$. Substituting these results in Eq. (21.28) we find that

$$\left(\frac{dE}{dt} \right)_{\text{scat}} = \frac{\sigma_{Tc}}{4\pi} \gamma^2 (4\pi U_R) \left(1 + \frac{1}{3} \frac{v^2}{c^2} \right) = \sigma_{Tc} \gamma^2 \left(1 + \frac{1}{3} \frac{v^2}{c^2} \right) U_R, \quad (21.32)$$

where we have used the relation $\langle E^2 \rangle = \langle B^2 \rangle = 4\pi U_R$. The incident radiation energy has been absorbed by the electron. The rate at which this happens is

$$\left(\frac{dE}{dt} \right)_{\text{abs}} = \sigma_{Tc} U_{\text{rad}}. \quad (21.33)$$

Hence, the net addition of energy to the photon field is:

$$\begin{aligned} P = \left(\frac{dE}{dt} \right) &= \left(\frac{dE}{dt} \right)_{\text{scat}} - \left(\frac{dE}{dt} \right)_{\text{abs}} = \sigma_{Tc} U_{\text{rad}} \left[\gamma^2 \left(1 + \frac{1}{3} \frac{v^2}{c^2} \right) - 1 \right] \\ &= \frac{4}{3} \sigma_{Tc} U_{\text{rad}} \gamma^2 \left(\frac{v}{c} \right)^2. \end{aligned} \quad (21.34)$$

It is remarkable how the various numerical factors play out correctly to give precisely the same coefficient $4/3$ in the final expression!

So whether we treat the thermal radiation field as a bunch of photons or as fluctuating electromagnetic fields, the final result is consistently the same.

All is well that ends well

Radiation reaction force, non-relativistic case

Appendix: The radiation reaction force g^i can be determined by using the criterion that the mean power radiated should be equal to the work done by the damping force. In the non-relativistic case, this leads to:

$$\left\langle \left(\frac{\Delta \mathcal{E}}{\Delta t} \right) \right\rangle = - \left\langle \left(\frac{2}{3} \right) q^2 a^2 \right\rangle = \langle \mathbf{f} \cdot \mathbf{v} \rangle, \quad (21.35)$$

when averaged over a period of time. Averaging a^2 over a time interval T , we get:

$$\begin{aligned}\langle a^2 \rangle &= \frac{1}{T} \int_0^T dt a^2 = \frac{1}{T} \int_0^T dt (\dot{\mathbf{v}} \cdot \dot{\mathbf{v}}) \\ &= \frac{1}{T} \int_0^T dt \left[\frac{d}{dt} (\mathbf{v} \cdot \dot{\mathbf{v}}) - \mathbf{v} \cdot \ddot{\mathbf{v}} \right] = \frac{1}{T} [\mathbf{v} \cdot \dot{\mathbf{v}}]_0^T - \langle \mathbf{v} \cdot \ddot{\mathbf{v}} \rangle. \quad (21.36)\end{aligned}$$

The first term vanishes as $T \rightarrow \infty$ for any bounded motion, giving $\langle a^2 \rangle = -\langle \mathbf{v} \cdot \ddot{\mathbf{v}} \rangle$. Using this, we see that $\mathbf{f}_{damp} = (2/3) q^2 \ddot{\mathbf{v}}$ in the non-relativistic case.

To obtain the corresponding relativistic expression for the four-force, we have to find a four-vector g^i which reduces to $(0, (2/3) q^2 \ddot{\mathbf{v}})$ in the rest frame of the charge. This condition is satisfied by any vector of the form $g^i = (2q^2/3) [(d^2 u^i/ds^2) + A u^i]$ where A is to be determined. To find A , we use the relativistic condition $g^i u_i = 0$ which should hold for any four-force. This gives $A = u^k (d^2 u_k/ds^2)$. Therefore:

*Radiation reaction
force, relativistic
case*

$$g^i = \left(\frac{2q^2}{3} \right) \left[\frac{d^2 u^i}{ds^2} + u^i u^k \frac{d^2 u_k}{ds^2} \right]. \quad (21.37)$$

The second term can be rewritten using

$$u^k \frac{da_k}{ds} = \frac{d}{ds} (u^k a_k) - a^k a_k = -a^k a_k, \quad (21.38)$$

since $u^k a_k = 0$. This gives another expression for g^i :

$$g^i = \frac{2}{3} q^2 \left[\frac{d^2 u^i}{ds^2} - u^i (a^k a_k) \right]. \quad (21.39)$$

Angular Momentum without Rotation

22

Electromagnetic fields exert forces on charged particles thereby changing the energy and momentum of the charged particles. If you now think of the charged particles and the electromagnetic field as making up a single closed system, then it follows that the energy and momentum supplied to the charged particle must have come from the electromagnetic field. This is indeed true and you must have learnt that the electromagnetic field possesses energy per unit volume (U) and momentum per unit volume (\mathbf{P}) given by

It sure has energy and momentum ...

$$U = \frac{1}{8\pi}(E^2 + B^2); \quad \mathbf{P} = \frac{1}{4\pi c}(\mathbf{E} \times \mathbf{B}). \quad (22.1)$$

However, what is not stressed adequately in text books is that electromagnetic fields — and even pretty simple ones — also possess angular momentum. Just as the electromagnetic field can exchange its energy and momentum with charged particles, it can also exchange its angular momentum with a system of charged particles, often leading to rather surprising results. In this chapter, we shall explore one such example.

... but angular momentum?!

One simple configuration [85] in which the exchange of angular momentum occurs is shown in Fig. 22.1. A plastic disk, located in the $x - y$ plane, is free to rotate about the vertical z -axis. A thin metallic ring of radius a , carrying a uniformly distributed charge Q , is embedded on the disk. Along the z -axis, there is a current-carrying solenoid producing a magnetic field \mathbf{B} contributing a total flux Φ . This initial configuration is completely static with a magnetic field \mathbf{B} confined within the solenoid and an electric field \mathbf{E} produced by the charge located on the ring. Let us suppose that the current source is disconnected, leading the magnetic field to die down. The change in the magnetic flux will lead to an electric field which will act tangential to the ring of charge, thereby giving it a torque. Once the magnetic field has died down, this torque would have resulted

A simple gadget

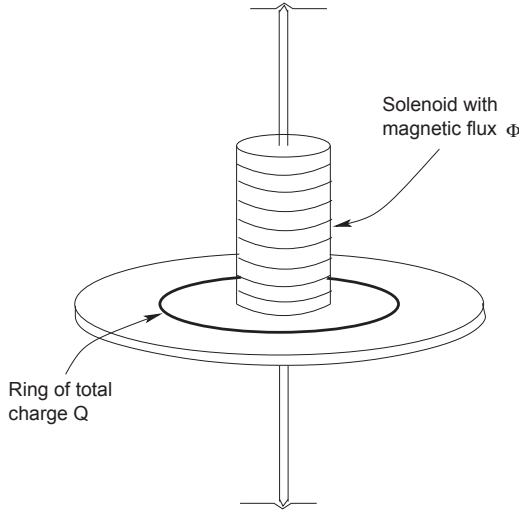


Fig. 22.1: A device to extract electromagnetic angular momentum by transferring it into rotational motion of charged particles. The circular disk with a ring of charge is free to rotate about the vertical axis. A coil of wire carrying current provides a solenoidal magnetic field near the axis in the vertical direction. Surprisingly, this static configuration (with the electric field of a charged ring and the magnetic field of the solenoid) stores certain amount of angular momentum. If the current in the solenoid is switched off, this angular momentum will be transferred to the ring of charge making the disk rotate.

in the disk spinning about the z -axis with a finite angular momentum. Where did the angular momentum come from?

The devil is in the details

It is quite obvious that the angular momentum in the initial field is what appears as the mechanical angular momentum of the rotating disk in the final stage. What is really interesting is to work this out and explicitly verify that the angular momentum is conserved (which Ref. [85] doesn't do!). I will now describe this calculation as well as some interesting issues which arise from it (There is large literature on this problem not all of which is illuminating; one place to start the search is from Refs. [86, 87]).

The final angular momentum is easy

The angular momentum of the final rotating disk can be computed easily. The rate of change of angular momentum $d\mathbf{L}/dt$ due to the torque acting on the ring of charge is along the z -axis, so we only need to compute its magnitude. This is given by:

$$\frac{dL}{dt} = aQE = \frac{Q}{2\pi} \oint \mathbf{E} \cdot d\mathbf{l} = -\frac{Q}{2\pi c} \frac{\partial \Phi}{\partial t}. \quad (22.2)$$

Here, \mathbf{E} is the tangential electric field generated due to the changing magnetic field and the last equality follows from Faraday's law. Integrating this equation and noting that the initial angular momentum of the disk and

the final magnetic flux are zero, we get

$$L = \frac{Q}{2\pi c} \Phi_{\text{initial}} . \quad (22.3)$$

It is interesting that the final angular momentum depends only on the total flux and not on other configurational details.

We now need to show that the initial *static* electromagnetic configuration had this much of stored angular momentum. We will first do this in a slightly unconventional manner and then indicate the connection to the more familiar approach.

Let us recall that the canonical momentum of a charge q located in a magnetic field is given by $\mathbf{p} - (q/c)\mathbf{A}$ where \mathbf{A} is the vector potential related to the magnetic field by $\mathbf{B} = \nabla \times \mathbf{A}$ and \mathbf{p} is the usual kinematic momentum. This suggests that one can associate with charges located in a magnetic field, a momentum $(q/c)\mathbf{A}$. For a distribution of charge, with a charge density ρ , the field momentum per unit volume will be $(1/c)\rho\mathbf{A}$. Hence, to a charge distribution located in a region of vector potential \mathbf{A} , we can attribute an angular momentum

Momentum in the presence of EM field

One way to define EM angular momentum

$$\mathbf{L}_A = \frac{1}{c} \int d^3\mathbf{x} \rho(\mathbf{x}) [\mathbf{x} \times \mathbf{A}(\mathbf{x})] . \quad (22.4)$$

In our problem, the charge distribution is confined to a ring of radius a with negligible magnetic field at the location of the charge. But the vector potential will exist outside the solenoid and the above expression can be non-zero. To compute this, let us use a cylindrical coordinate system with (r, θ, z) as the coordinates. We will choose a gauge in which the vector potential has only the tangential component; that is, only A_θ is non-zero. Using

\mathbf{A} exist where \mathbf{B} doesn't!

Choose a gauge

$$\oint \mathbf{A} \cdot d\mathbf{l} = \Phi , \quad (22.5)$$

where Φ is the total magnetic flux, we get $2\pi r A_\theta = \Phi$ for a line integral of A around any circle. Hence $A_\theta = \Phi/(2\pi r)$. This can be written in a nice vectorial form as

$$\mathbf{A} = \frac{\Phi}{2\pi r^2} (\hat{\mathbf{z}} \times \mathbf{r}) , \quad (22.6)$$

where $\hat{\mathbf{z}}$ is the unit vector in the z -direction. When we substitute this expression in Eq. (22.4) and calculate the angular momentum, the integral gets contribution only from a circle of radius a . Further, using the identity, $\mathbf{r} \times (\hat{\mathbf{z}} \times \mathbf{r}) = \hat{\mathbf{z}} r^2$, we get the result that

$$\mathbf{L}_A = \frac{Q}{2\pi c} \Phi_{\text{initial}} \hat{\mathbf{z}} , \quad (22.7)$$

which is exactly the final angular momentum which we computed in Eq. (22.3). Rather nice!

The initial angular momentum, as expected

The good news ...

However, the above elementary derivation, as well as the expression for electromagnetic angular momentum in Eq. (22.4), raises several intriguing issues. On the positive side, it makes the vector potential a very tangible quantity, something which we learnt from relativity and quantum mechanics but could never be clearly demonstrated within the context of classical electromagnetism. In the process, it also gives a physical meaning to the field momentum $(q/c)\mathbf{A}$ which is somewhat mysterious in conventional approaches. On the flip side, one should note that \mathbf{A} , by the very definition, is gauge dependent and one would have preferred a definition of the electromagnetic angular momentum which is properly gauge invariant.

... and the bad news

A more conventional definition

It is, of course, possible to write down another, more conventional, expression for the electromagnetic angular momentum. Given the density of electromagnetic momentum, \mathbf{P} , we can define the corresponding angular momentum density as $\mathbf{x} \times \mathbf{P}$. Integrating it over all space should give the angular momentum associated with the electromagnetic field. Since the momentum density \mathbf{P} involves only the electric and magnetic fields, the resulting expressions are automatically gauge invariant. This leads to a definition of angular momentum given by

$$\mathbf{L}_{\text{EM}} = \frac{1}{4\pi c} \int d^3\mathbf{x} [\mathbf{x} \times (\mathbf{E} \times \mathbf{B})], \quad (22.8)$$

which just replaces the momentum density $\rho\mathbf{A}/c$ in Eq. (22.4) by $(\mathbf{E} \times \mathbf{B}/4\pi c)$. It is easy to verify that, as momentum *densities*, these two expressions are *unequal* in general. But what is relevant, as far as our computation goes, is the integral over the whole space of these two expressions. If these two expressions differ by terms which vanish when integrated over the whole space, then we have an equivalent gauge invariant definition of field angular momentum.

Do the two definitions give the same result?

It turns out that this is indeed the case in any static configuration if we choose to describe the magnetic field in a gauge with $\nabla \cdot \mathbf{A} = 0$. One can then show that

$$\frac{1}{4\pi}(\mathbf{E} \times \mathbf{B})^\alpha = \frac{1}{4\pi}(\mathbf{E} \times (\nabla \times \mathbf{A}))^\alpha = \rho\mathbf{A}^\alpha + \frac{\partial V^{\beta\alpha}}{\partial x^\beta}, \quad (22.9)$$

where $V^{\beta\alpha}$ is a complicated second rank tensor built out of field variables. While one can provide a proof of Eq. (22.9) using vector identities (you should try it out!), it is a lot faster and neater to use four dimensional notation and special relativity to get this result.

The proof, in relativistic notation

We begin with the expression for the momentum density of the electromagnetic field in terms of the stress tensor T^{ab} of the electromagnetic field. The T^{00} component of this tensor is proportional to the energy density of the electromagnetic field while the $T^{0\alpha}$ component is proportional

to the momentum density P^α . More precisely,

$$T_0^\alpha = -\frac{1}{4\pi}(\mathbf{E} \times \mathbf{B})^\alpha = -cP^\alpha. \quad (22.10)$$

On the other hand, the electromagnetic stress tensor can be written in terms of the four dimensional field tensor F^{ab} in the form $T_0^\alpha = (1/4\pi)F^{\alpha\beta}F_{0\beta}$. We now manipulate this expression using the facts that (i) the configuration is static and (ii) the vector potential satisfies the gauge condition $\nabla \cdot \mathbf{A} = \partial_\alpha A^\alpha = 0$, to prove Eq. (22.9). Using the definition of the field tensor in terms of the four vector potential, $F_{ij} = \partial_i A_j - \partial_j A_i$, we can write:

$$\begin{aligned} T_0^\alpha &= \frac{1}{4\pi}F^{\alpha\beta}F_{0\beta} = \frac{1}{4\pi}(\partial^\alpha A^\beta - \partial^\beta A^\alpha)F_{0\beta} \\ &= \frac{1}{4\pi}(\partial^\alpha A^\beta)F_{0\beta} - \frac{1}{4\pi}\partial^\beta(F_{0\beta}A^\alpha) + A^\alpha \frac{\partial^\beta F_{0\beta}}{4\pi} \\ &= \frac{1}{4\pi}(-\partial^\alpha A^\beta \partial_\beta A_0) - \frac{1}{4\pi}\partial^\beta(F_{0\beta}A^\alpha) + A^\alpha \frac{\partial^\beta F_{0\beta}}{4\pi}. \end{aligned} \quad (22.11)$$

To arrive at the second line, we have performed an integration by parts and to obtain the third line, we have used $\partial_0 A_\beta = 0$ since the configuration is time independent. We next use the result $\partial^\beta F_{0\beta} = -\nabla \cdot \mathbf{E} = -4\pi\rho$ in the last term and perform another integration by parts in the first term, using the gauge condition $\nabla \cdot \mathbf{A} = \partial_\alpha A^\alpha = 0$. This gives

$$T_0^\alpha = -\rho A^\alpha - \frac{1}{4\pi}\partial_\beta[A_0\partial^\alpha A^\beta - A^\alpha\partial^\beta A_0]. \quad (22.12)$$

We thus find that:

$$cP^\alpha = \rho A^\alpha + \partial_\beta V^{\beta\alpha}; \quad V^{\beta\alpha} \equiv \frac{1}{4\pi}[A_0\partial^\alpha A^\beta - A^\alpha\partial^\beta A_0], \quad (22.13)$$

which proves the equivalence between the two expressions for electromagnetic momentum density ($c\mathbf{P}$ and $\rho\mathbf{A}$), when used in integrals over all space, provided the second term vanishes sufficiently fast. For the case we are discussing, this is indeed true.

Final result

From the result in Eq. (22.9), it is easy to see that, in our example, we will get the same result irrespective of whether we use \mathbf{L}_A or \mathbf{L}_{EM} . This is because, when we integrate the expressions in Eq. (22.9) over all space, the term involving $V^{\beta\alpha}$ can be converted to a surface term at infinity which does not contribute.

The first observation of what we now call Brownian motion was probably made by the Dutch physicist Jan Ingenhauze, the discoverer of photosynthesis. In 1785, he put alcohol to good use by sprinkling powdered charcoal on it and observing it under a microscope. The name Brownian motion for the random perambulation of the particles comes from Robert Brown, who published an extensive investigation of similar phenomena in 1828. This was eventually heralded as evidence for molecular nature of matter, and figured crucially in the award of the 1926 Nobel Prize in physics to Jean Perrin for determining the Avogadro number.

Molecules: stand up and be counted!

The term “random walk”, on the other hand, appears to have been first coined by Carl Pearson in 1905, the same year in which Einstein published his paper on Brownian motion. Pearson was interested in providing a simple model for the spread of mosquito infestation in a forest — which goes to show, right at the outset, the generality of the process! Pearson’s letter to *Nature* was answered by Lord Rayleigh who had solved this problem earlier in the case of sound waves in heterogeneous material. Independent of all this, Louis Bachelor was developing the theory of random walks in his remarkable doctoral thesis *La theorie de la speculation* published in 1900. Here, the random walk was suggested as a model for a financial time series, which has, until recently, helped physicists get Wall Street jobs (with the consequences we all now know only too well!). This brief glimpse of history already shows the occurrence of the random walk in widely different contexts. (An entertaining discussion of history is available in Refs. [88, 89]; also see Ref. [90].)

What is common to: the spread of mosquitoes, sound waves and the flow of money?

Let us begin by reviewing the simplest of all random walks in which a particle moves from the origin, taking steps of length ℓ , with each step being in a random direction uncorrelated with the previous one. The displacement of the particle after N steps is given by

$$\mathbf{x} = \sum_{n=1}^N \mathbf{x}_n, \quad (23.1)$$

where

$$|\mathbf{x}_n| = \ell; \quad \langle \mathbf{x}_n \rangle = 0; \quad \langle \mathbf{x}_n \cdot \mathbf{x}_m \rangle = \ell^2 \delta_{nm} . \quad (23.2)$$

The first equation shows that each step has a constant magnitude. The second and third equations denote averaging over a probability distribution by the symbol $\langle \dots \rangle$ and quantifies the uncorrelated nature of the directions of the steps. From these, we can immediately obtain two key results of such a random walk. First, $\langle \mathbf{x} \rangle = 0$. Further, we have

$$\sigma^2 \equiv \langle \mathbf{x}^2 \rangle = \left\langle \left(\sum_{n=1}^N \mathbf{x}_n \right)^2 \right\rangle = \sum_{n,m=1}^{\infty} \langle \mathbf{x}_n \cdot \mathbf{x}_m \rangle = N \ell^2 . \quad (23.3)$$

This shows that the key characteristic of the random walk, viz., the root-mean-square displacement σ grows as \sqrt{N} .

*No limits for
 dx/\sqrt{dt}*

We can think of ℓ as Δx , denoting the magnitude of the displacement between any two consecutive steps. If the time interval between the steps is Δt , then $\sigma \propto \sqrt{N}$ suggests that $(\Delta x)^2/\Delta t$ remains constant in the continuum limit. Clearly, the random walk corresponds to a curve without a definite slope in the continuum limit and, in fact, the continuum limit needs to be taken with some care [91]. This is one of the many reasons why random walks are fascinating.

To see how such a continuum limit emerges in this context, we should generalize the concept of random walk slightly by assuming that the probability for the particle to take a step given by the vector $\Delta \mathbf{y}$ is given by some function $p(\Delta \mathbf{y})$ with the properties

$$\begin{aligned} \langle \Delta y^i \rangle &\equiv \int d^D \Delta \mathbf{y} [\Delta y^i p(\Delta \mathbf{y})] = 0; \\ \langle \Delta y^i \Delta y^j \rangle &\equiv \int d^D \Delta \mathbf{y} [\Delta y^i \Delta y^j p(\Delta \mathbf{y})] = \langle (\Delta y)^2 \rangle \frac{\delta^{ij}}{D} . \end{aligned} \quad (23.4)$$

where $i, j, \dots = 1, 2, \dots, D$ denotes the components of the vector. Let $P_N(\mathbf{x})$ be the probability that the net displacement is \mathbf{x} after N steps. Then, since the steps are uncorrelated, we have the elementary relation:

$$P_N(\mathbf{x}) = \int d^D \Delta \mathbf{y} P_{N-1}(\mathbf{x} - \Delta \mathbf{y}) p(\Delta \mathbf{y}) . \quad (23.5)$$

To obtain the continuum limit, we will assume that a Taylor series expansion of $P_{N-1}(\mathbf{x} - \Delta \mathbf{y})$ is possible so that we can write (assuming summation over repeated indices):

$$\begin{aligned}
P_N(\mathbf{x}) &\cong \int d^D \Delta y p(\Delta \mathbf{y}) \left\{ P_{N-1}(\mathbf{x}) - \Delta y^i \partial_i P_{N-1}(\mathbf{x}) \right. \\
&\quad \left. + \frac{1}{2} \Delta y^i \Delta y^j \partial_i \partial_j P_{N-1}(\mathbf{x}) \right\} \\
&= P_{N-1}(\mathbf{x}) + \frac{\langle (\Delta y)^2 \rangle}{2D} \nabla^2 P_{N-1}(\mathbf{x}) , \tag{23.6}
\end{aligned}$$

where we have used Eq. (23.4). In the continuum limit, we will denote the total time which has elapsed since the beginning of random walk by $t = N\Delta t$ and define a continuum probability density by $\rho(\mathbf{x}, t) = \rho(\mathbf{x}, N\Delta t) \equiv P_N(\mathbf{x})$. Since we can take $(\partial\rho/\partial t)$ as the limit $[P_N(\mathbf{x}) - P_{N-1}(\mathbf{x})]/\Delta t$ when $\Delta t \rightarrow 0$, we get from Eq. (23.6), the result:

$$\frac{\partial \rho}{\partial t} = K \nabla^2 \rho , \tag{23.7}$$

where we have defined a ('diffusion') coefficient $K \equiv \langle (\Delta y)^2 \rangle / 2D\Delta t$. The continuum limit exists if we can treat K as a constant when $\Delta t \rightarrow 0$. This, clearly, is equivalent to $\langle (\Delta y)^2 \rangle / \Delta t$ being finite in the continuum limit as we indicated earlier. This is quite different from the usual continuum limits we are accustomed to in physics in which the ratio of the differentials of *the same order* are replaced by a derivative. This warns us that something non-trivial is going on.

The trick ...

Note that the final equation we have obtained, of course, is the diffusion equation which can also be written as $(\partial\rho/\partial t) = -\nabla \cdot \mathbf{J}$ where the current $\mathbf{J} = -K\nabla\rho$ arises due to the gradient in the particle density. (In this form, we can even consider a situation with spatially varying diffusion coefficient K .) This indicates that diffusive processes in physics can be modeled at the microscopic level by a random walk of the discrete constituent element. The diffusion equation is also unique in the sense that it is not invariant under time reversal; diffusion gives you a direction of time — which is another remarkable feature that arises in the continuum limit.

... leading to the diffusion equation

The diffusion equation, Eq. (23.7), being a linear equation, can be solved by Fourier transforming both sides. Denoting the Fourier transform of $\rho(\mathbf{x}, t)$ by $\rho(\mathbf{k}, t)$, it is easy to show that $\rho(\mathbf{k}, t) = \exp(-Kk^2 t)$. Taking a Fourier transform, we get the fundamental solution to the diffusion equation (which is essentially the Green's function) to be

$$\rho(\mathbf{x}, t) = \frac{e^{-x^2/4Kt}}{(4\pi Kt)^{D/2}} . \tag{23.8}$$

This shows how particles located close to the origin at $t = 0$ spread out in the course of time. The mean square spread is clearly proportional to Kt which is the residue of the discrete result $\sigma^2 \propto N$.

*Application:
Diffusion in
velocity space*

The diffusion of a particle need not always take place in the real 3-dimensional space. An interesting phenomenon which occurs in plasmas as well as gravitating systems — wherein the long range, inverse square forces act between particles — involves diffusion in the *velocity space*. A simple version of this can be described as follows.

Consider a near homogeneous distribution of gravitationally interacting particles (e.g., stars in a globular cluster). When two stars scatter off each other with an impact parameter b , each one undergoes a typical acceleration Gm/b^2 acting for a time b/v . In any one such scattering, a typical star will acquire a “kick” in the velocity space of magnitude $(\delta v_\perp \approx Gm/bv)$, $\delta v_\perp \ll v$. The effect of a large number of such collisions is to make the star perform a random walk in the velocity space. The net mean-square-velocity induced by collisions with impact parameters in the range $(b, b+db)$ in a time interval Δt is the product of the mean number of scatterings in time (Δt) and $(\delta v_\perp)^2$. The former is given by the number of scatterers in the volume $(2\pi b db)(v\Delta t)$. Therefore,

$$\langle (\delta v_\perp)^2 \rangle = (2\pi b db) (v\Delta t) n \left(\frac{Gm}{bv} \right)^2. \quad (23.9)$$

where n is the number density of scatterers. The total mean-square transverse velocity due to all stars is found by integrating over b within some range (b_1, b_2) :

$$\begin{aligned} \langle (\delta v_\perp)^2 \rangle_{\text{total}} &\simeq \Delta t \int_{b_1}^{b_2} (2\pi b db) (vn) \left(\frac{G^2 m^2}{b^2 v^2} \right) \\ &= \frac{2\pi n G^2 m^2}{v} \Delta t \ln \left(\frac{b_2}{b_1} \right). \end{aligned} \quad (23.10)$$

*All divergences
arise from incom-
plete physics*

We again see the signature of the random walk in $\langle \delta v_\perp^2 \rangle \propto \Delta t$. The logarithmic factor shows that we cannot take $b_1 = 0, b_2 = \infty$, and need to use some physical criteria to fix b_1 and b_2 . It is reasonable to take $b_2 \simeq R$, the size of the system; as regards b_1 , notice that the velocity changes per collision can become comparable to v itself when $b \simeq b_c \simeq (Gm/v^2)$ and our diffusion approximation breaks down. It is, therefore, reasonable to take $b_1 \simeq b_c \simeq (Gm/v^2)$. Then $(b_2/b_1) \simeq (Rv^2/Gm) = N (Rv^2/GM) \simeq N$ for a system in virial equilibrium. From Eq. (23.10), we see that this effect is important over time-scales (Δt) which are long enough to make $\langle (\delta v_1)^2 \rangle_{\text{total}} \simeq v^2$. Using this condition and solving for (Δt) , we get:

$$(\Delta t) \simeq \frac{v^3}{2\pi G^2 m^2 n \ln N}. \quad (23.11)$$

This is the timescale for gravitational relaxation in such systems (or electromagnetic relaxation in plasmas) and the $\ln N$ factor arises due to diffusion in velocity space.

The entire process can be described by a diffusion equation in velocity space — or so it would seem at first sight. Further thought, however, shows that if we describe the process by a diffusion equation in velocity space, it will make the root-mean-square velocities of *every* particle in the system increase as \sqrt{t} as time goes on, which violates some sacred notions in physics. (For a description of the curious history behind these discoveries, see e.g. Ref. [92].) This is one key difference between diffusion in real space, compared to velocity space, and there must exist another process which prevents this.

You don't want this to go on and on

This process is called dynamical friction. To understand this process, consider a particle (“star”) which moves with a speed V that is significantly larger than the root mean square speed of the cloud of stars around it. In the rest frame of the fast star, on the average, other stars will stream past it and will be deflected towards it. This will produce a slight density enhancement of stars behind the fast star. This density enhancement produces the necessary extra force to reduce the speed V of the star. This dynamical friction ensures that no runaway disaster occurs in the velocity space.

The other side: dynamical friction

If we take both the processes into account, the evolution in the velocity space is described by an equation which is a variant of what is called the Fokker-Planck equation. We can describe the diffusion in the velocity (or, equivalently, momentum) space, that obeys standard conservation laws, through a source term which is a divergence of a current in the momentum space. Hence, the evolution of the distribution function will be governed by an equation of the form

$$\frac{df}{dt} = \frac{\partial f}{\partial t} + \mathbf{v} \cdot \frac{\partial f}{\partial \mathbf{x}} - \nabla \phi \cdot \frac{\partial f}{\partial \mathbf{v}} = - \frac{\partial J^\alpha}{\partial p^\alpha} . \quad (23.12)$$

The form of the current J_α can be determined by considering the elementary collisional process, and one obtains [102, 103] the result

$$J_\alpha(\ell) = \frac{B_0}{2} \int d\ell' \left\{ f \frac{\partial f'}{\partial \ell_\beta} - f' \frac{\partial f}{\partial \ell_\beta} \right\} \cdot \left\{ \frac{\delta_{\alpha\beta}}{k} - \frac{k_\alpha k_\beta}{k^3} \right\}; \quad \mathbf{k} = \ell - \ell' , \quad (23.13)$$

where

$$B_0 = 4\pi G^2 m^5 L; \quad L = \int_{b_1}^{b_2} \frac{db}{b} = \ln \left(\frac{b_2}{b_1} \right) \approx \ln N . \quad (23.14)$$

In this current in Eq. (23.13), the term proportional to f leads to dynamical friction while the term proportional to $\partial f / \partial \ell_\beta$ leads to the increase in the velocity dispersion. The form in Eq. (23.13) is quite elegant and, by

inspection, we can conclude that the current vanishes for the Maxwellian distribution which should arise as the steady state configuration. I will not bother to derive the above equation for you (if you are interested, look through the references in Box 23.1) but will illustrate the nature of this equation using a simpler one.

A simplified version of this equation, which contains the essential features for our purpose, is given by

$$\frac{\partial f(v,t)}{\partial t} = \frac{\partial}{\partial v} \left\{ \frac{\sigma^2}{2} \frac{\partial f}{\partial v} + (\alpha v) f \right\}. \quad (23.15)$$

Fokker-Planck in a simple case

The first term on the right hand side has the standard form of a diffusion current proportional to the gradient in the velocity space. As time goes on, this term will cause the mean square velocities of particles to increase in proportion to t , inducing the ‘random walk’ in the velocity space. Under the effect of this term, *all* the particles in the system will have their $\langle v^2 \rangle$ increasing without bound. This unphysical situation is avoided by the presence of the second term $(\alpha v f)$ which describes the dynamical friction. The combined effect of the two terms is to drive f to a Maxwellian distribution with an effective temperature $(k_B T) = (\sigma^2/\alpha)$ and $(\partial f/\partial t) = 0$. In such a Maxwellian distribution, the gain made in (Δv^2) due to diffusion is exactly balanced by the losses due to dynamical friction. When two particles scatter, one gains the energy lost by the other; on the average, we may say that the one which has lost the energy has undergone dynamical friction while the one which gained energy has achieved diffusion to higher v^2 . The cumulative effect of such phenomena is described by the two terms in Eq. (23.15).

The solution

The above features can be illustrated by explicitly solving Eq. (23.15). Suppose we take an initial distribution $f(v,0) = \delta(v - v_0)$ peaked at a velocity v_0 . The solution of Eq. (23.15) with this initial condition is:

$$f(v,t) = \left[\frac{\alpha}{\pi \sigma^2 (1 - e^{-2\alpha t})} \right]^{1/2} \exp \left[-\frac{\alpha (v - v_0 e^{-\alpha t})^2}{\sigma^2 (1 - e^{-2\alpha t})} \right], \quad (23.16)$$

which is a Gaussian with the mean $\langle v \rangle = v_0 e^{-\alpha t}$ and dispersion $\langle v^2 \rangle - \langle v \rangle^2 = (\sigma^2/\alpha)(1 - e^{-2\alpha t})$. At late times ($t \rightarrow \infty$), the mean velocity $\langle v \rangle$ goes to zero while the velocity dispersion becomes (σ^2/α) . Thus, the equilibrium configuration is a Maxwellian distribution of velocities with this particular dispersion, for which the right hand side of Eq. (23.15) vanishes.

To see the effect of the two terms individually on the initial distribution $f(v,0) = \delta(v - v_0)$, we can set α or σ to zero. When $\alpha = 0$, we get pure diffusion:

$$f_{\alpha=0}(v,t) = \left(\frac{1}{2\pi\sigma^2 t} \right)^{1/2} \exp \left\{ -\frac{(v - v_0)^2}{2\sigma^2 t} \right\}. \quad (23.17)$$

Nothing happens to the steady velocity v_0 ; but the velocity dispersion increases in proportion to t representing a random walk in the velocity space. If, on the other hand, we set $\sigma = 0$, then we get

$$f_{\sigma=0}(v, t) = \delta(v - v_0 e^{-\alpha t}). \quad (23.18)$$

Now there is no spreading in velocity space (no diffusion); instead the friction steadily decreases $\langle v \rangle$.

Box 23.1: History: Landau's derivation of dynamical friction

One of the key results in Chandrasekhar's book [93] is the derivation of the collisional relaxation time. He essentially obtains the result in Eq. (23.11) after devoting about 25 pages (from pages 48 to 73) for the algebraic derivation which includes a 'three-dimensional' picture! For comparison, the same result had been obtained earlier by James Jeans [94] in 1929 using about 3 pages (pages 317 to 320) in his book. The result was doubtless known to many others and — in fact — the explicit use of $\ln N$ in the timescale for soft collisions exists in a 1938 paper of Ambartsumian [95]. Chandrasekhar defends his elaborate calculation of this previously known result by saying: "Though the physical ideas were correctly formulated by Jeans a completely rigorous evaluation of the time of relaxation was not available until recently". Chandrasekhar does not seem to have been bothered by the fact that any estimation of time of relaxation will necessarily be uncertain by factors of order unity both because of the variation of density — Chandrasekhar assumes a constant density star cluster — and by the uncertainties in the upper and lower cut-offs inside the logarithm.

There is, however, a more interesting twist to this tale which illustrates one of the rare occasions in which Chandrasekhar completely missed a key physical effect. As mentioned in the text, if all the stars in a globular cluster continue random walking in velocity space, it will violate some sacred principles of physics. It seems that Chandrasekhar realized this soon after — but only after — the publication of his book on stellar dynamics. He addresses this issue and obtains the expression for dynamical friction as a *separate* physical phenomenon in his works published shortly afterwards [96–98].

Curiously enough, the elegant expression in Eq. (23.13), giving both the dynamical friction and diffusion at one go, was already known before Chandrasekhar's work! These results were first obtained and published — in 1936, about six years before Chandrasekhar's work was published — by Landau [99]. (He was dis-

cussing Coulomb interactions in a plasma but everything can be trivially translated to gravitational interaction.)

Strangely, the elegance and power of this result was not appreciated, occasionally even by plasma physicists. A detailed discussion of this approach [100] by Rosenbluth et al in 1957 cites Chandrasekhar's work but not Landau's though they have a citation to Cohen et al. [101] with a comment "A more complete list of references is given in this..." The paper by Cohen et al. does cite Landau's paper but it is clear they have not understood the result at all, because they say that, in Landau's work, "... the important terms representing dynamical friction which should appear in the diffusion equation are set equal to zero as a result of certain approximations". which is, of course, incorrect. Landau, in the usual elegant but terse style, has captured all the essential physics. (A textbook derivation of this result in the context of plasmas [102] as well as gravitating systems [103] is now available.)

*A more general
random walk ...*

Returning to the discrete case, we can make another useful generalization of Eq. (23.5) by assuming that $p(\Delta \mathbf{y})$ itself depends on N so that the fundamental equation becomes

$$P_N(\mathbf{x}) = \int d^D y P_{N-1}(\mathbf{x} - \Delta \mathbf{y}) p_N(\Delta \mathbf{y}) . \quad (23.19)$$

This equation, which is a convolution integral, is also easy to solve in Fourier space in which the convolution integral becomes a product. If we denote by $P_N(\mathbf{k})$ and $p_N(\mathbf{k})$ the Fourier transforms of $P_N(\mathbf{x})$ and $p_N(\Delta \mathbf{y})$ respectively, then this equation becomes $P_N(\mathbf{k}) = P_{N-1}(\mathbf{k}) p_N(\mathbf{k})$. Iterating this N times and normalizing the initial probability by assuming the particle was at the origin, we get:

$$P_N(\mathbf{k}) = \prod_{n=1}^N p_n(\mathbf{k}) . \quad (23.20)$$

... which is solvable

Performing an inverse Fourier transform, we find the solution to our problem to be

$$P_N(\mathbf{x}) = \int \frac{d^D k}{(2\pi)^D} e^{i\mathbf{k} \cdot \mathbf{x}} \prod_{n=1}^N p_n(\mathbf{k}) . \quad (23.21)$$

Again, it is possible to make some general comments if the individual probability distributions $p_n(\Delta \mathbf{y})$ satisfy some reasonable conditions. Consider, for simplicity, that $p_n(\Delta \mathbf{y})$ is peaked at the origin and dies down smoothly and monotonically for large $|\Delta \mathbf{y}|$. Then, its Fourier transform will also be peaked around the origin in k -space and will die down for large values of $|\mathbf{k}|$. Further, because the probability is normalized, we have

the condition $p_n(\mathbf{k} = 0) = 1$. When we take a product of N such functions, the resulting function will again have the value unity at the origin. But as we go away from the origin, we are taking the product of N numbers each of which is less than unity. Clearly, when $N \rightarrow \infty$, the product of $p_n(\mathbf{k})$ will have significant support only close to the origin.

A simple trick ...

The non-trivial assumption we will now make is that $p_n(\mathbf{k})$ has a smooth curvature at the origin of the Fourier space and is not ‘cuspy’. Then, near the origin in Fourier space, we can approximate

$$p_n(\mathbf{k}) \simeq 1 - \frac{1}{2} \alpha_n^2 k^2 \simeq e^{-(1/2) \alpha_n^2 k^2} , \quad (23.22)$$

with some constant α_n . The product then becomes:

$$\prod_{n=1}^N p_n(\mathbf{k}) = \exp -\frac{1}{2} k^2 \sum_{n=1}^N \alpha_n^2 \equiv \exp -\frac{N}{2} \sigma^2 k^2 , \quad (23.23)$$

where we have defined

$$\sigma^2 = \frac{1}{N} \sum_{n=1}^N \alpha_n^2 . \quad (23.24)$$

In this limit, the final Fourier transform in Eq. (23.21) will give a Gaussian in \mathbf{x} with $\langle x^2 \rangle \propto N$.

... to prove the central limit theorem!

An observant reader would have noticed that we have essentially proved a variant of the central limit theorem for the sum $(\mathbf{x}_1 + \mathbf{x}_2 + \dots \mathbf{x}_N)$ of N independently distributed random variables, each having its own probability distribution $p_n(\mathbf{x}_n)$. In fact, the joint probability for these variables to be in some given interval is given by the product, $p_n(\mathbf{x}_n) d^D \mathbf{x}_n$ over all $n = 1, 2, \dots, N$. The probability for their sum to be \mathbf{x} is given by

$$P_N(\mathbf{x}) = \int \prod_{n=1}^N p_n(\mathbf{x}_n) d^D \mathbf{x}_n \delta_D(\mathbf{x} - \sum \mathbf{x}_n) , \quad (23.25)$$

where the Dirac delta function ensures that the sum of the random variables is \mathbf{x} . We write the Dirac delta function in Fourier space to obtain

$$\begin{aligned} P_N(\mathbf{x}) &= \int \frac{d^D k}{(2\pi)^D} e^{i\mathbf{k} \cdot \mathbf{x}} \prod_{n=1}^N \int d^D \mathbf{x}_n p_n(\mathbf{x}_n) e^{-i\mathbf{k} \cdot \mathbf{x}_n} \\ &= \int \frac{d^D k}{(2\pi)^D} e^{i\mathbf{k} \cdot \mathbf{x}} \prod_{n=1}^N p_n(\mathbf{k}) , \end{aligned} \quad (23.26)$$

which is identical to the result we obtained earlier in Eq. (23.21).

A classic example in which our analysis (and the central limit theorem) fails is given by the case in which each of the probability distributions $p_n(\Delta \mathbf{y})$ is given by a Lorentzian

When the central limit theorem fails

$$p_n(\Delta \mathbf{y}) = \frac{(\beta/\pi)}{(\Delta \mathbf{y})^2 + \beta^2} . \quad (23.27)$$

The Fourier transform now gives $p_n(\mathbf{k}) = \exp(-\beta|\mathbf{k}|)$. Clearly the approximation in Eq. (23.22) fails for this function, since it is ‘cuspy’ due to a linear term in $|\mathbf{k}|$ near the origin. We can, of course, carry out the analysis in Eq. (23.21) to get

$$P_N(\mathbf{x}) = \int \frac{d^D k}{(2\pi)^D} e^{i\mathbf{k} \cdot \mathbf{x}} e^{-N\beta|\mathbf{k}|} = \frac{(N\beta/\pi)}{|\mathbf{x}|^2 + (N^2\beta^2)} . \quad (23.28)$$

More is not different

We have the result that the probability distribution for the final displacement is identical to the probability distribution of individual steps when the latter is a Lorentzian — except for the (expected) scaling of the width.

The main reason for the central limit theorem to fail in this case is that the Lorentzian distribution has a diverging second moment. You should remember this the next time you think of the full width at half maximum of a Lorentzian as “similar to” the width of a Gaussian! There are physical situations, (e.g., one called anomalous diffusion), which can be modeled along these lines. They are characterized by random walks in which every once in a while the particle takes a large step because of the slow decrease in the probability $p(\Delta\mathbf{y})$.

Random walk on a lattice

Quite often, one also considers the random walk on a lattice of specific shape, the simplest being the D-dimensional cube. Here, the particle hops from one site of the lattice to a nearby site along any one of the axes with the lattice spacing taken to be unity for simplicity. In this case, the Fourier integrals in Eq. (23.21) become Fourier series, and we get:

$$P_N(\mathbf{x}) = \int_{-\pi}^{\pi} \frac{d^D k}{(2\pi)^D} [\cos(\mathbf{k} \cdot \mathbf{x})] \prod_{n=1}^N p_n(\mathbf{k}) , \quad (23.29)$$

where all the integrals are in the range $(-\pi, \pi)$ and \mathbf{x} is a vector with integer valued components. If $p_n(\mathbf{k})$ is independent of n , and the hops in all directions from any site are equally likely, then $p(\mathbf{k}) = (1/D)(\cos k_1 + \cos k_2 + \cdots \cos k_D)$ and we get:

$$P_N(\mathbf{x}) = \int_{-\pi}^{\pi} \frac{d^D k}{(2\pi)^D} [\cos(\mathbf{k} \cdot \mathbf{x})] \left(\frac{1}{D} \sum_{j=1}^D \cos k_j \right)^N . \quad (23.30)$$

As a cross check, we can reproduce the standard result for the one dimensional lattice using Eq. (23.30). In this case $x = J$, with J being a positive or negative integer. After N steps when the particle has taken n_L steps to the left of origin and n_R steps to the right, we have $n_L + n_R = N$ and $n_R - n_L = J$. Solving this, we get $n_R = (1/2)(N + J)$, $n_L = (1/2)(N - J)$. The probability that out of N steps n_L was to the left and n_R was to the right is the same as getting, say, n_L heads while tossing N coins, and is

given by:

$$P_N(J) = \frac{1}{2^N} {}^N C_{n_L} = \frac{1}{2^N} \frac{N!}{((1/2)(N+J))!((1/2)(N-J))!} . \quad (23.31)$$

You can amuse yourself by proving that the same expression is also given by the integral in Eq. (23.30) for $D = 1$,

$$P_N(J) = \int_{-\pi}^{\pi} \frac{dk_1}{(2\pi)} [\cos(k_1 J)] (\cos k_1)^N , \quad (23.32)$$

as it should. The result in Eq. (23.30) will be useful in the next chapter when we address some interesting dimension-dependent properties of random walks (and an unexpected connection with electrical networks!).

More on Random Walks: Circuits and a Tired Drunkard

24

The general formula for the probability $P_N(\mathbf{x})$ for a particle to be found at position \mathbf{x} after N steps, obtained in the last chapter [see Eq. (23.30)], depends on the dimension of space D in which the random walk takes place. (It also depends on the geometry of the lattice, but for simplicity, we will only consider cubic lattice in D -dimensions.) Do the crucial features of random walks depend on the dimension D ? At first sight, one might have thought that the random walk in, say, $D = 1, 2, 3$ will behave in essentially the same manner. Curiously enough, this is not the case!

Surprise, surprise:
 $3 \neq 2 \neq 1$

The dimensional dependence of the the random walk can be illustrated [104] by studying the phenomenon known as *recurrence*. Recurrence refers to the probability for the the random walking particle to come back to the origin — where it started from — in the course of its perambulation, when we wait for infinite time. Let u_n denote the probability that a particle returns to the origin on the n th step and let \mathcal{R} be the expected number of times it returns to the origin. Clearly,

$$\mathcal{R} = \sum_{n=0}^{\infty} u_n . \quad (24.1)$$

We can now distinguish between two different scenarios. If the series in Eq. (24.1) diverges, then the mean number of returns to the origin is infinite and we say that the the random walk is recurrent. If the series is convergent, leading to a finite \mathcal{R} , then we say that the the random walk is transient.

This idea is reinforced by the following alternative interpretation of \mathcal{R} . Suppose u is the probability for the particle to return to the origin. Then, the normalized probability for it to return exactly k times is $u^k(1-u)$. The mean number of returns to the origin is, therefore,

Another perspective

$$\mathcal{R} = \sum_{k=1}^{\infty} k u^{k-1} (1-u) = (1-u)^{-1} . \quad (24.2)$$

Obviously, if $\mathcal{R} = \infty$, then $u = 1$, showing that the the random walker will definitely return to the origin. But if $\mathcal{R} < \infty$, then $u < 1$ and it is not certain that the particle will ever come back home.

Let us compute u_n and \mathcal{R} for random walks in $D = 1, 2, 3$ dimensions with the lattice spacing set to unity for simplicity. From Eq. (23.30), setting $\mathbf{x} = 0$, we have:

$$u_n = \int_{-\pi}^{\pi} \frac{d^D k}{(2\pi)^D} \left(\frac{1}{D} \sum_{j=1}^D \cos k_j \right)^n. \quad (24.3)$$

Doing the sum in Eq. (24.1) we get

$$\mathcal{R} = \sum_{n=0}^{\infty} u_n = \int_{-\pi}^{\pi} \frac{d^D k}{(2\pi)^D} \left(1 - \frac{1}{D} \sum_{j=1}^D \cos k_j \right)^{-1}. \quad (24.4)$$

We want to ascertain whether this integral is finite or divergent. Clearly, the divergence, if any, can only arise due to its behaviour near the origin in k -space. Using the Taylor series expansion of the cosine function, we see that, near the origin, we have the behaviour:

$$\mathcal{R} \approx 2D \int_{k \approx 0} \frac{dk_1 dk_2 \dots dk_D}{(2\pi)^D} (k_1^2 + k_2^2 \dots k_D^2)^{-1} \propto \frac{2D}{(2\pi)^D} \int_{k \approx 0} \frac{k^{D-1} dk}{k^2}. \quad (24.5)$$

*A drunken man
will definitely come
home, in the long
run, but a drunken
bird may or may
not!*

The dimension dependence is now obvious. In $D = 1, 2$ the integral is divergent and $\mathcal{R} = \infty$; so we conclude that the random walk in $D = 1, 2$ is recurrent and the particle will definitely return to the origin if it walks forever. But in $D = 3$, \mathcal{R} is finite and the walk is non-recurrent. There is finite probability that the particle will come back to the origin but there is also a finite probability that it will not.

The mean number of recurrences in $D = 3$ is given by — what is known as — the Watson integral

$$\mathcal{R} = \frac{3}{(2\pi)^3} \int_{-\pi}^{\pi} dk_1 \int_{-\pi}^{\pi} dk_2 \int_{-\pi}^{\pi} dk_3 [3 - (\cos k_1 + \cos k_2 + \cos k_3)]^{-1}, \quad (24.6)$$

which is notoriously difficult to evaluate analytically. Since the answer happens to be

$$\mathcal{R} = \frac{\sqrt{6}}{32\pi^3} \Gamma\left(\frac{1}{24}\right) \Gamma\left(\frac{5}{24}\right) \Gamma\left(\frac{7}{24}\right) \Gamma\left(\frac{11}{24}\right), \quad (24.7)$$

you anyway need to look it up in a table so one might as well do the integral numerically (which is trivial in *Mathematica*, say) and get $\mathcal{R} \approx 1.5164$, giving the return probability $u \approx 0.3405$. This integral was

first evaluated by Watson [105] in terms of elliptic integrals and a “simpler” result was obtained by Glasser and Zucker later on [106].

In the case of $D = 1$ or 2 , it is also easy to obtain u_n explicitly by using a combinatorics argument. In 1-dimension, the particle can return to the origin only if it has taken an even number of steps, half to the right and half to the left. The probability for this is clearly

D = 1 is easy

$$u_{2n} = {}^{2n}C_n \frac{1}{2^{2n}} . \quad (24.8)$$

For sufficiently large n , we can use Stirling’s approximation for factorials ($n! \approx \sqrt{2\pi n} e^{-n} n^n$) to get $u_{2n} \approx 1/\sqrt{\pi n}$. The series in Eq. (24.1) involves the asymptotic sum which is divergent:

$$m = \sum_n u_{2n} \approx \sum_n \frac{1}{\sqrt{\pi n}} = \infty . \quad (24.9)$$

Obviously, the 1-dimensional random walk is recurrent.

Interestingly, the result for $D = 2$ turns out to be just the square of the result for $D = 1$. The integral in Eq. (24.3) becomes, for $D = 2$:

D = 2 is just the square of D = 1!

$$u_n(\mathbf{x}) = \frac{1}{(2\pi)^D} \frac{1}{2^n} \int_{-\pi}^{\pi} dk_1 \int_{-\pi}^{\pi} dk_2 (\cos k_1 + \cos k_2)^n . \quad (24.10)$$

If you now change the variables of integration to $(k_1 + k_2)$ and $(k_1 - k_2)$, it is easy to show that this integral becomes the product of the two integrals for the $D = 1$ case, giving

$$u_{2n} = \left[\frac{1}{2^{2n}} {}^{2n}C_n \right]^2 , \quad (24.11)$$

which is the square of the result for $D = 1$. Now, the series in Eq. (24.1) will be dominated asymptotically by

$$\mathcal{R} \approx \sum_n \frac{1}{\pi n} = \infty , \quad (24.12)$$

thus making the $D = 2$ random walk recurrent. You might guess at this stage that in 3- D , the asymptotic series will involve a sum over $n^{-3/2}$ (and hence will converge) making the 3- D random walk non-recurrent. This is partially true and the 3-dimensional series is *bounded from above* by the sum over $n^{-3/2}$. But the 3-dimensional case is *not* the product of three 1-dimensional cases.

No, it doesn’t work for D = 3!

We now turn our attention to another curious result. Summing $P_N(\mathbf{x})$ over all N , one can construct the quantity $P(\mathbf{x})$ which is the net probability of reaching a location \mathbf{x} . Using Eq. (23.30) and doing the geometric sum,

we find this quantity, in $D = 2$, to be:

$$P(\mathbf{x}) = \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} \frac{dk_1 dk_2}{(2\pi)^2} [\cos(\mathbf{k} \cdot \mathbf{x})] \left(1 - \frac{1}{2}(\cos k_1 + \cos k_2) \right)^{-1}. \quad (24.13)$$

Consider now the expression

$$\begin{aligned} R &= \frac{1}{2}(P(\mathbf{x}) - P(\mathbf{0})) \\ &= \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} \frac{dk_1 dk_2}{8\pi^2} \frac{[1 - \cos(\mathbf{k} \cdot \mathbf{x})]}{[1 - (\cos k_1 + \cos k_2)/2]}. \end{aligned} \quad (24.14)$$

Incredibly enough, this provides the solution to a completely different problem! Consider a grid of 1 ohm resistors connected between the lattice sites of an infinite, two-dimensional square lattice. It turns out that R is the effective resistance between the lattice point \mathbf{x} and the origin. Let us see how this comes about by analyzing the grid of resistors.

*The real pleasure
of doing physics;
what you thought is
different is the same!*

Let a node \mathbf{x} in the infinite planar square lattice be denoted by two integers (m, n) and let a current $I_{m,n}$ be injected at that node. The flow of current will induce a voltage at each node and, using Kirchoff's and Ohm's laws for the 1 ohm resistors we can write the relation:

$$\begin{aligned} I_{m,n} &= (V_{m,n} - V_{m+1,n}) + (V_{m,n} - V_{m-1,n}) + (V_{m,n} - V_{m,n+1}) + (V_{m,n} - V_{m,n-1}) \\ &= 4V_{m,n} - V_{m+1,n} - V_{m-1,n} - V_{m,n+1} - V_{m,n-1}, \end{aligned} \quad (24.15)$$

where $V_{m,n}$ is the potential at the node (m, n) due to the current. This equation can again be solved by introducing the Fourier transform on the discrete lattice. If we write

$$I_{m,n} = \frac{1}{4\pi^2} \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} dk_1 dk_2 I(k_1, k_2) e^{i(mk_1 + nk_2)} \quad (24.16)$$

$$V_{m,n} = \frac{1}{4\pi^2} \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} dk_1 dk_2 V(k_1, k_2) e^{i(mk_1 + nk_2)}, \quad (24.17)$$

then one can obtain from Eq. (24.15) the result in the Fourier space:

$$I(k_1, k_2) = 2V(k_1, k_2) [2 - \cos(k_1) - \cos(k_2)]. \quad (24.18)$$

Suppose a current of 1 amp is injected at $(0,0)$, and (-1) amp at (N, M) . Then $I_{m,n} = \delta_{m,n} - \delta_{m-M, n-N}$, leading to

$$I(k_1, k_2) = 1 - e^{-i(Mk_1 + Nk_2)}, \quad (24.19)$$

so that Eq. (24.18) gives the voltage to be

$$V(k_1, k_2) = \frac{1}{2} \frac{1 - e^{-i(Mk_1 + Nk_2)}}{2 - \cos(k_1) - \cos(k_2)}. \quad (24.20)$$

The equivalent resistance between nodes (0,0) and (M,N) with the a flow of unit current is just the voltage difference between the nodes:

The answer

$$\begin{aligned}
 R_{M,N} &= V_{0,0} - V_{M,N} \\
 &= \frac{1}{4\pi^2} \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} dk_1 dk_2 V(k_1, k_2) \left[1 - e^{i(Mk_1 + Nk_2)} \right] \\
 &= \frac{1}{4\pi^2} \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} dk_1 dk_2 \frac{1}{2} \frac{(1 - e^{-i(Mk_1 + Nk_2)})(1 - e^{i(Mk_1 + Nk_2)})}{2 - \cos(k_1) - \cos(k_2)} \\
 &= \frac{1}{4\pi^2} \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} dk_1 dk_2 \frac{1 - \cos(Mk_1 + Nk_2)}{2 - \cos(k_1) - \cos(k_2)}, \quad (24.21)
 \end{aligned}$$

which is exactly the same as the integral in Eq. (24.14)!

The infinite grid of square lattice resistors is a classic problem and the effective resistance between two adjacent nodes is a “trick question” that is a favourite of examiners. The answer [0.5 ohm] can be found by trivial superposition but the effective resistance between arbitrary nodes cannot be obtained by such tricks. In fact, the effective resistance between two diagonal nodes of the basic square — the (0,0) and (1,1), say — is given by the integral

*Clever tricks are fine
but hard work wins
every time!*

$$\begin{aligned}
 R_{1,1} &= \frac{1}{4\pi^2} \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} dk_1 dk_2 \frac{1 - \cos(k_1 + k_2)}{2 - \cos(k_1) - \cos(k_2)} \\
 &= \frac{1}{\pi^2} \int_0^{\pi} dk_1 \int_0^{\pi} dk_2 \frac{1 - \cos(k_1) \cos(k_2)}{2 - \cos(k_1) - \cos(k_2)}. \quad (24.22)
 \end{aligned}$$

The second equality is obtained by noting that the denominator is an even function and hence only the even part of $\cos(k_1 + k_2)$ needs to be kept in the numerator. Once the entire integral is an even function, we can change the limits to 0 and π and multiply by 4. The resulting integral is fairly straightforward but a bit tedious and can be done as follows. You split it as two integrals and use the standard results:

$$\int_0^{\pi} dv \frac{1}{[2 - \cos(u)] - \cos(v)} = \frac{\pi}{\sqrt{[2 - \cos(u)]^2 - 1}}, \quad (24.23)$$

and

$$\int_0^{\pi} dv \frac{\cos(v)}{[2 - \cos(u)] - \cos(v)} = \pi \left[\frac{2 - \cos(u)}{\sqrt{[2 - \cos(u)]^2 - 1}} - 1 \right], \quad (24.24)$$

to evaluate the integral over, say, k_2 . After some simplification, this reduces the integral to the form:

$$R_{1,1} = \frac{1}{\pi} \int_0^{\pi} dk_1 \frac{[1 - \cos(k_1)]^2}{\sqrt{[2 - \cos(k_1)]^2 - 1}}. \quad (24.25)$$

Substituting $x = 1 - \cos k_1$, this integral reduces to:

$$R_{1,1} = \frac{1}{\pi} \int_0^2 dx \frac{x}{\sqrt{4-x^2}} = \frac{2}{\pi} . \quad (24.26)$$

Clearly, the equivalent resistance between two diagonal lattice points of the infinite grid is a transcendental number involving π . (Next time someone lectures you on the power of clever arguments, ask her to get Eq. (24.26) by clever arguments!)

But why does this work ? What is the correspondence between the random walk on a lattice and resistor networks ? There are different levels of sophistication at which one can answer this question. There is a large literature on this subject and an entire book [107] dealing with this subject exists. The mathematical reason has to do with the fact that both the random walk probability to visit a node and the voltage on a node (which does not have any current injected or removed) are harmonic functions. These are functions whose value at any given node is given by the average of the value of the function on the adjacent lattice sites. This is obvious in the case of the random walk because a particle which reaches the node (m, n) must have hopped to that node with equal probability from one of the neighbouring nodes $(m \pm 1, n \pm 1)$. In the case of the resistor network, the same result is obtained from Eq. (24.15) when $I_{mn} = 0$. If you now inject the voltages 1 and 0 at two specific nodes A and B , then the voltage at any other node X can be interpreted as the probability that a random walker starting at X will get to A before B . One can then use this interpretation to make a formal connection between voltage distribution in an electric network and a random walk problem. The interested reader can find more in the book [107] referred above.

We now go back to the random walk in the continuum for which we had obtained the result in the last chapter, which, specialized to one dimension, is given by:

$$P_N(x) = \int_{-\infty}^{\infty} \frac{dk}{(2\pi)} e^{ikx} \prod_{n=1}^N p_n(k) . \quad (24.27)$$

Surely, any drunkard will get tired as he walks?

We now consider a situation in which the steps are random and uncorrelated but their lengths are decreasing monotonically. In particular, we will assume that each step length is a fraction λ of the previous one, with $\lambda < 1$, and the first step is of unit length. It is clear that $P_N(x)$ is now given by

$$P_N(x) = \int_{-\infty}^{\infty} \frac{dk}{(2\pi)} e^{ikx} \prod_{n=1}^N \cos(k\lambda^n) . \quad (24.28)$$

We can now study the limit of $N \rightarrow \infty$ and ask how the probability $P_{\infty}(x, \lambda) \equiv P_{\lambda}(x)$ (with a slight change in notation) is distributed. (This interesting topic does not seem to have been explored in sufficient detail. A good discussion is available in ref. [108, 109].)

This probability distribution has very beautiful and unexpected features. To begin with, when λ is less than $(1/2)$, the support of the function $P_\lambda(x)$ (i.e., the range of x for which $P_\lambda(x)$ is non-zero) is a Cantor set! On the other hand, when $(1/2) \leq \lambda < 1$, there is a countably infinite set of λ values for which $P_\lambda(x)$ is singular, with it being smooth for almost all other values of λ . The most interesting case occurs when λ takes the value of the golden ratio, $\lambda = g \equiv (\sqrt{5} - 1)/2$. This $P_g(x)$ is riddled with singularities but shows a remarkable self-similar behaviour. I will now describe some of these features.

Cantor set ?! Where does that spring from?

Let us first consider some cases for which one can obtain simple analytic results. Take, for example, the case of $\lambda = 1/2$ which is on the borderline between the two behaviours. In this case, the relevant infinite product is given by:

$$\prod_{n=1}^{\infty} \cos \frac{k}{2^n} = \frac{\sin k}{k} . \quad (24.29)$$

(This is a cute result which you can prove as follows: Write

Maths trick

$$\cos \frac{k}{2^n} = \frac{1}{2} \frac{\sin(k/2^{n-1})}{\sin(k/2^n)} , \quad (24.30)$$

take a product of N terms canceling out the sines and then take the limit $N \rightarrow \infty$.) Since the Fourier transform of $(\sin k/k)$ is just a uniform distribution, we get the tantalizing result that $P(x)$ is just a uniform distribution in the interval $(-1, 1)$ and zero elsewhere !

Similar methods also work for $\lambda = 2^{-1/2}, 2^{-1/4} \dots$ etc. For example, when $\lambda = 2^{-1/2}$, the infinite product is

$$\prod_{n=1}^{\infty} \cos \frac{k}{2^{n/2}} = \left(\frac{\sin k}{k} \right) \left(\frac{\sin \sqrt{2}k}{\sqrt{2}k} \right) . \quad (24.31)$$

The Fourier transform of this involves a convolution of two rectangular distributions and is easily seen to be a triangular probability distribution. For the case of $\lambda = 2^{-1/m}$, the relevant product again can be evaluated in a similar manner and the distribution will have continuous derivatives up to order $(m - 1)$ while the m th derivative will be discontinuous at $2m$ points. Clearly as $m \rightarrow \infty$, the distribution becomes more and more smooth and approaches the Gaussian limit of the standard random walk.

There is a clever way of understanding the end point distribution for random walks in which the step length varies as 2^{-n} or 3^{-n} . In the first case, the final resting place for our tired drunkard is

Another insight into the result

$$S = \sum_{n=1}^{\infty} a_n 2^{-n} , \quad (24.32)$$

where a_n is a random variable taking the values ± 1 with equal probability. From this, it is easy to see that

$$S + 1 = \sum_{n=1}^{\infty} a_n 2^{-n} + \sum_{n=1}^{\infty} 2^{-n}. \quad (24.33)$$

So

$$\frac{1}{2}(S + 1) = \sum_{n=1}^{\infty} \frac{1}{2}(a_n + 1)2^{-n} = \sum_{n=1}^{\infty} \omega_n 2^{-n}, \quad (24.34)$$

where ω_n takes the values 0 or 1 with equal probability. We now notice that this is just the expression for a number in the interval $[0, 1]$ in base 2 with ω_n denoting the digits 0 or 1 in the binary expansion. Hence the probability distribution for $(1/2)(S + 1)$ is uniform in the interval $[0, 1]$. It follows that the probability distribution for S is uniform in the interval $[-1, +1]$.

A similar trick works when the step size falls as 3^{-n} . In this case, we have the relation:

$$S = \sum_{n=1}^{\infty} a_n 3^{-n}, \quad a_n \in \{-1, 1\}; \quad S + \frac{1}{2} = \sum_{n=1}^{\infty} t_n 3^{-n}, \quad t_n \in \{0, 2\}. \quad (24.35)$$

Cantor set from base-3

We now see that $S + (1/2)$ is given by the representation of a number in the interval $[0, 1]$ written in base 3 but has *only* the digits 0 and 2 appearing in it. This is actually the Cantor set. Therefore, S is distributed over the Cantor set constructed from the interval $[-(1/2), (1/2)]$ by removing the middle term.

Let us next try to understand why we get something as strange as a Cantor set when $\lambda < 1/2$. One way of doing this is as follows. We first note that one can think of the geometric random walk as a random *map* given by the equation

$$x' = \pm 1 + \lambda x, \quad (24.36)$$

which describes how the position of the particle changes in a single step. This is obvious if you substitute x' for x on the right hand side and iterate. You will find that the map is equivalent, after infinite steps, to the random sum

$$x = \sum_n \varepsilon_n \lambda^n; \quad \varepsilon_n = \pm 1. \quad (24.37)$$

A useful recursion

Further, we note that our random walk problem satisfies a simple recursion relation. If $P_\lambda(x, N)$ is the probability to be at location x after N steps, then it is obvious that

$$P_\lambda(x, N) = \frac{1}{2} \left[P_\lambda \left(\frac{x-1}{\lambda}, N-1 \right) + P_\lambda \left(\frac{x+1}{\lambda}, N-1 \right) \right]. \quad (24.38)$$

Since everything converges for $\lambda < 1/2$, we can take the limit of $N \rightarrow \infty$ in this equation to obtain

$$P_\lambda(x) = \frac{1}{2} \left[P_\lambda \left(\frac{x-1}{\lambda} \right) + P_\lambda \left(\frac{x+1}{\lambda} \right) \right]. \quad (24.39)$$

If we now define the probability measure $M_\lambda(a, b)$ for x to be found in the interval (a, b) by the integral:

$$M_\lambda(a, b) = \int_a^b dx P_\lambda(x), \quad (24.40)$$

then we get the corresponding recursion relation to be:

$$2M_\lambda(a, b) = M_\lambda \left(\frac{a-1}{\lambda}, \frac{b-1}{\lambda} \right) + M_\lambda \left(\frac{a+1}{\lambda}, \frac{b+1}{\lambda} \right). \quad (24.41)$$

Using $M_\lambda(a, b)$ has the advantage that it smoothens out the singularities in $P_\lambda(x)$. It is obvious from Eq. (24.41) that the support of M_λ lies in the interval $[-x_{\max}, x_{\max}]$ with $x_{\max} = (1 - \lambda)^{-1}$. When $\lambda < 1/2$, our map in Eq. (24.36) transforms this interval to the union of two non-overlapping intervals given by

$$\left[-\frac{1}{(1-\lambda)}, \frac{(1-2\lambda)}{(1-\lambda)} \right], \left[-\frac{1-2\lambda}{(1-\lambda)}, \frac{1}{(1-\lambda)} \right]. \quad (24.42)$$

If we use the map in Eq. (24.36) again to either of these sub-intervals, they, in turn, get mapped into further non-overlapping sub-intervals. If we continue these iterations an infinite number of times, we obtain the final support for M_λ which is clearly a Cantor set!

Iterate to Cantor set

A more intuitive interpretation of this bifurcation can be provided along the following lines. Suppose the first step in the random walk is to the right. So, the maximum displacement of the subsequent walk is $\lambda/(1-\lambda)$. Therefore the end point of the walk must necessarily lie in the region

$$\left[1 - \frac{\lambda}{1-\lambda}, 1 + \frac{\lambda}{1-\lambda} \right]. \quad (24.43)$$

We note that the left edge of this region is positive when $\lambda < 1/2$; so the support of $P_\lambda(x)$ has got divided into two non-overlapping regions just after one step. Clearly, the same kind of bifurcation occurs at each step finally leading to a Cantor set.

What about the singular behaviour which arises when $\lambda > 1/2$? This result, in contrast, is extraordinarily hard to analyse. But one can qualitatively see why singular behaviour might arise for certain special values of λ (which, by no means, is exhaustive). Consider the subset of λ values

Difficult, not completely solved

which satisfies the equation

$$1 - \sum_{n=1}^N \lambda^n = 0, \quad (24.44)$$

which can be viewed as a random walk with the first step of unit length to the right, followed by N steps to the left, such that we end up at the origin. By solving the polynomial equation for $N = 2, 3, 4, \dots$, we get the values

$$\lambda = \left\{ \frac{1}{2}(\sqrt{5} - 1), 0.544, 0.519, \dots \right\}, \quad (24.45)$$

where the first entry is the inverse of the golden ratio $g \cong 0.618$. This positional degeneracy of returning to the origin (in which the points are reached by different random walks with same number of steps) is the basic reason for the singular behaviour of $P_\lambda(x)$.

The Golden Walk

As we said before, the largest of these values $\lambda = g$, which is the inverse of the golden ratio, has very special properties. It has a self-similar structure because of which the probability distribution in the interval $J^0 \equiv [-g, g]$ reproduces the full distribution if we rescale the length by a factor g^{-3} and the probability by a factor 3. It turns out that these results arise because this probability distribution has an infinite number of symmetries underlying such a distribution but this is way too complicated for us to discuss here.

Gravitational Instability of the Isothermal Sphere

25

A generic problem related to the establishment of thermodynamic equilibrium can be stated as follows: Consider a large number (N) of particles, interacting through a two-body potential $U(\mathbf{x} - \mathbf{y})$ and confined in a region of volume V . We start off the particles with a generic set of initial positions and velocities and let them interact (“collide”) with each other as well as with the boundary of the volume V . We are interested in the very late time behaviour of such a system. In particular, we are often interested in the kind of equilibrium configuration to which such a system might evolve into at sufficiently late times.

The general problem

The result will clearly depend on the nature of the interaction, specified by $U(\mathbf{x} - \mathbf{y})$ as well as the other parameters. If $U(\mathbf{x} - \mathbf{y})$ is a short range potential representing intermolecular forces and if E is sufficiently high, then the system will relax towards a Maxwellian distribution of velocities and a nearly uniform density in space. The velocity distribution will have characteristic temperature $T \simeq 2E/3N$ and we are assuming that this T is higher than the ‘boiling point’ of the ‘liquid’ made of these particles. If not, the eventual equilibrium state will be a mixture of matter in the liquid and vapour state. (Note that we use units with $k_B = 1$ throughout.) All this is part of standard lore in statistical mechanics.

Standard phases

What happens if $U(\mathbf{x} - \mathbf{y})$ is due to gravitational interaction of the particles? What are the different phases in which matter can exist in such a case? In this chapter, we will discuss some of the peculiar effects that arise in this context. Let us begin by recalling some details of the standard statistical mechanics applied to systems with short range interactions.

In the study of laboratory systems involving short range interaction between constituent particles, a central quantity which we use is the entropy functional $S(E, V)$ that gives the entropy of the system as a function of energy and volume. This, in turn, is related to the density of states of the system $g(E)$ by $S(E) = \ln g(E)$ with

$$g(E) \equiv \frac{d\Gamma(E)}{dE}; \quad \Gamma(E) \equiv \int dp dq \theta[E - H(p, q)], \quad (25.1)$$

where $H(p, q)$ is the Hamiltonian and $\theta(z)$ is the Heaviside function with $\theta(z) = 1$ for $z \geq 0$ and zero otherwise. (We will suppress exhibiting the explicit dependence of various quantities on the volume V when it is not relevant.) In this microcanonical description of the system, the temperature and the pressure can be obtained by

$$T(E) = \left(\frac{\partial S}{\partial E} \right)^{-1}; \quad P = T \left(\frac{\partial S}{\partial V} \right), \quad (25.2)$$

which shows that the relation between temperature and energy can be determined once we know the Hamiltonian $H(p, q)$ of the system. For example, an ideal gas of N particles with $H \propto \sum p_i^2$ will lead to the familiar relations

$$\Gamma \sim V^N E^{3N/2} \sim g(E); \quad T(E) = (2E/3N); \quad P/T = N/V, \quad (25.3)$$

when $N \gg 1$.

Quite often one uses the equivalent canonical description based on the partition function $Z(T)$ given by the Laplace transform of the density of states

$$Z(T) = \int dE g(E) \exp[-\beta E] = \int dp dq \exp[-\beta H(p, q)], \quad (25.4)$$

where $\beta = 1/T$. In this case, one determines the (mean) energy and pressure by the relations

$$\bar{E} = -(\partial \ln Z / \partial \beta); \quad \bar{P} = T(\partial \ln Z / \partial V). \quad (25.5)$$

For systems which obey extensivity of energy, (viz., when total energy of the system is the sum of its parts to a high degree of accuracy) the canonical and microcanonical descriptions will lead to the same physical quantities to the accuracy $\mathcal{O}(\ln N/N)$ where N is the number of degrees of freedom of the system.

Let us now consider what happens in self-gravitating systems [110]. The first casualty is the equivalence between canonical and microcanonical descriptions which fails for systems with gravitational interaction mainly because energy is not an extensive parameter for such systems (see e.g., [111]). If a large gravitating system is divided into two parts the total energy cannot be expressed as the sum of the energies of the two parts; the (gravitational) interaction energy between the parts of the system makes a significant contribution to the total energy due to the long range nature of gravity. Hence the fundamental description of gravitating systems has to be based on microcanonical ensemble and any use of canonical ensemble (in some occasions) needs to be justified by specific physical considerations.

Ideal gas: quick recap

canonical \neq microcanonical!

This inequivalence of the two ensembles should also be obvious from the fact that systems in canonical ensemble cannot exhibit negative specific heat while self-gravitating systems often do. The first result is obvious from Eq. (21.2) which tells us that the specific heat C_V for a system in canonical ensemble is given by

In canonical ensemble $C_V > 0$, while gravity makes $C_V < 0$

$$C_V \equiv \frac{\partial \bar{E}}{\partial T} = -\beta^2 \frac{\partial \bar{E}}{\partial \beta} = \beta^2 (\Delta E)^2 > 0, \quad (25.6)$$

The second result follows from the fact that, for gravitating systems in steady state, virial theorem gives $2K + U = 0$, where K is the kinetic energy and U is the potential energy. This implies $E = -K$ where $E = K + U$ is the total energy. Since the temperature is proportional to the kinetic energy of random motion K , it follows that gravitating systems in steady state, obeying virial theorem, have negative specific heat. Obviously, one needs to be careful in using standard results from statistical mechanics of laboratory systems to describe gravitating systems.

The sensible — though not always practical — thing to do is to use the most basic of the ensembles, viz. the microcanonical ensemble to describe the gravitating systems. To do this, we need to evaluate the density of states in Eq. (25.1). This integral will diverge in the absence of two relevant cut-offs. First is the cut-off at large distances which is required to confine high energy particles from moving to large distances. This, of course, is not special to self-gravitating systems; even an ideal gas of particles will have a divergent density of states if it is not confined by a box of volume V . The second cut-off is at short distances to prevent particles from approaching each other arbitrarily closely thereby releasing large amount of gravitational potential energy, $-Gm^2/r$, as $r \rightarrow 0$. Once again, such a situation arises even in the case of plasmas in which quantum mechanical considerations will provide an effective short distance cut-off. For gravitating systems relevant to astrophysics there is usually some other physical process, say, arising from the finite size of the self-gravitating objects, which will provide this cut-off.

Two physical cut-offs

Given a large distance cut-off R and short distance cut-off a one can, in principle, compute the density of states and the thermodynamic behaviour of such a system. The two cut-offs define two natural energy scales $E_1 = -Gm^2/a$ and $E_2 = -Gm^2/R$ with $a \ll R$. On the other hand, the application of virial theorem to such a system will lead to a relation of the form

$$2K + U = 3PV + U_0, \quad (25.7)$$

where K is the kinetic energy of the particles, U is the gravitational potential energy due to standard ($-r^{-1}$) potential, P is the pressure exerted by the particles on the confining volume and U_0 is the correction to the virial due to the short distance cut-off. Broadly speaking, the different phases of

Many phases of gravitating systems

the gravitating systems can be related [111] to the different ways in which this condition is satisfied:

(a) When the energy of the system is such that $E \gg E_2$, gravity is irrelevant and the system behaves like a gas confined by a container. In this high temperature phase with positive specific heat Eq. (25.7) is satisfied with $2K \approx 3PV$ and the other two terms are sub-dominant; i.e., $U \ll K$ and $U_0 \ll 3PV$.

(b) When $E_1 \ll E \ll E_2$, the system is unaffected either by the confining box or by the short distance cut-off. In this phase with *negative* specific heat, it is dominated entirely by gravity and Eq. (25.7) is satisfied by $2K + U \approx 0$ with the other two terms being sub-dominant ($U_0 \ll U$, $3PV \ll U$). Since canonical ensemble cannot lead to negative specific heat, the description in canonical and microcanonical ensembles differ drastically in this regime. In canonical ensemble, the negative specific heat region is replaced by a rapid phase transition.

(c) As we approach lower energies ($E \rightarrow E_1$) the hard core nature of the particles begins to be felt and the gravity is resisted by the other physical processes. This will lead to a low temperature hard core condensate in which Eq. (25.7) is satisfied by $U \approx U_0$ with the other two terms being sub-dominant ($2K \ll U$, $3PV \ll U_0$).

The existence of negative specific heat phase is characteristic of the inherent instability of self-gravitating system. As the system evolves, it has a tendency to form a centrally condensed core with $U \approx U_0$ releasing large amount of energy which puts the remaining part of the system into a high temperature phase that will exist as a halo around the core. The particles in this halo will be bouncing off the walls of the container in the form of a high temperature gas with a cold core existing as a centrally condensed body.

The above description should convince you that the statistical mechanics of self gravitating systems is quite complex and it is not easy to make analytical progress with microcanonical ensemble. The next best thing is to use an approximation called the mean-field theory. I will now describe this approach in the context of self-gravitating systems.

Physical kinetics

Consider a system described by a distribution function $f(\mathbf{x}, \mathbf{p}, t)$ such that $\int d^3\mathbf{x} d^3\mathbf{p}$ denotes the total mass in a small phase space volume. We assume that the evolution of the distribution function is given by some equation (usually called the Boltzmann equation) of the form $df/dt = C(f)$. The term $C(f)$ on the right hand side describes the effect of collisions between the particles in the system. While the precise form of $C(f)$ may be complicated, it is usually assumed that the collisional evolution of f , driven by $C(f)$, satisfies two reasonable conditions: (a) The total mass and energy of the system are conserved and (b) the mean field entropy,

defined by,

$$S = - \int f \ln f d^3 \mathbf{x} d^3 \mathbf{p} , \quad (25.8)$$

does not decrease (and in general increases). If you are unfamiliar with this expression, here is a recap: In the standard derivation of the Boltzmann distribution, one extremises the function $S = - \sum n_i \ln n_i$ of the occupation numbers n_i , subject to the constraints, on total energy and number. In the continuum limit one works with f rather than n_i and the summation over i becomes an integral over the phase space, leading to Eq. (25.8). For any such system, we can obtain the equilibrium form of f by extremising the entropy while keeping the total energy and mass constant using two Lagrange multipliers. This is a standard exercise in statistical mechanics and the resulting distribution function is the usual Boltzmann distribution governed by:

$$f(\mathbf{x}, \mathbf{v}) \propto \exp \left[-\beta \left(\frac{1}{2} v^2 + \phi \right) \right]; \quad \phi(\mathbf{x}) = \int d^3 \mathbf{y} U(\mathbf{x}, \mathbf{y}) \rho(\mathbf{y}) . \quad (25.9)$$

Integrating over velocities, we get the closed system of integral equations for the density distribution:

The equilibrium density distribution

$$\rho(\mathbf{x}) = \int d^3 \mathbf{v} f = A \exp(-\beta \phi(\mathbf{x})); \quad \phi(\mathbf{x}) = \int d^3 \mathbf{y} U(\mathbf{x}, \mathbf{y}) \rho(\mathbf{y}) . \quad (25.10)$$

The final result is quite understandable: It is just the Boltzmann factor for the density distribution: $\rho \propto \exp(-\beta \phi)$ where ϕ is the potential energy at a given location due to the distribution of particles. One could have almost written this down “by inspection”!

The description so far is independent of the nature of the potential U (except for one important caveat which we will discuss right at the end). In the case of gravitational interaction, Eq. (25.10) becomes:

$$\rho(\mathbf{x}) = A \exp(-\beta \phi(\mathbf{x})); \quad \phi(\mathbf{x}) = -G \int \frac{\rho(\mathbf{y}) d^3 \mathbf{y}}{|\mathbf{x} - \mathbf{y}|} . \quad (25.11)$$

The integral equation (25.11) for $\rho(\mathbf{x})$ can be easily converted to a differential equation for $\phi(\mathbf{x})$ by taking the Laplacian of the second equation — leading to $\nabla^2 \phi = 4\pi G \rho$ — and using the first. We then get $\nabla^2 \phi \propto e^{-\beta \phi}$. If we now consider the spherically symmetric case, this reduces to:

Differential is easier than Integral, as equations go

$$\nabla^2 \phi = \frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\phi}{dr} \right) = 4\pi G \rho_c e^{-\beta[\phi(r) - \phi(0)]} , \quad (25.12)$$

called the isothermal sphere equation. (One can actually prove that among all solutions to Eq. (25.11), the spherically symmetric one extremises the S in Eq. (25.8).) The constants β and ρ_c (the central density) have to be

fixed in terms of the total number (or mass) of the particles and the total energy. Given the solution to this equation, which represents an extremum of the entropy, all other quantities can be determined. As we shall see, this system shows several peculiarities.

*Get rid of
inessentials*

To analyse Eq. (25.12), it is convenient to introduce length, mass and energy scales by the definitions

$$L_0 \equiv (4\pi G \rho_c \beta)^{1/2}, \quad M_0 = 4\pi \rho_c L_0^3, \quad \phi_0 \equiv \beta^{-1} = \frac{GM_0}{L_0}. \quad (25.13)$$

All other physical variables can be expressed in terms of the dimensionless quantities $x \equiv (r/L_0)$, $n \equiv (\rho/\rho_c)$, $m = (M(r)/M_0)$, $y \equiv \beta[\phi - \phi(0)]$ where $M(r)$ is the mass inside a sphere of radius r . These variables satisfy the easily derived equations:

$$y' = m/x^2; \quad m' = nx^2; \quad n' = -mn/x^2, \quad (25.14)$$

where the prime denotes the derivative with respect to x . In terms of $y(x)$, the isothermal equation, Eq. (25.12), becomes

$$\frac{1}{x^2} \frac{d}{dx} \left(x^2 \frac{dy}{dx} \right) = e^{-y}, \quad (25.15)$$

with the boundary condition $y(0) = y'(0) = 0$.

*One solution, but
not what we want*

Let us consider the nature of solutions to this equation. By direct substitution, we see that $n = (2/x^2)$, $m = 2x$, $y = 2 \ln x$ satisfy Eq. (25.14) and Eq. (25.15). This simple solution, however, is singular at the origin and hence is not physically admissible. The importance of this solution lies in the fact that – as we will see – all other (physically admissible) solutions tend to this solution [111, 112] for large values of x . This asymptotic behavior of all solutions shows that the density decreases as $(1/r^2)$ for large r implying that the mass contained inside a sphere of radius r increases as $M(r) \propto r$ at large r . Of course, in our case, the system is enclosed in a spherical box of radius R with a given mass M .

*A useful trick
to know*

To find non-singular solutions that satisfy the boundary conditions $y(0) = y'(0) = 0$, we first note that Eq. (25.15) is invariant under the transformation $y \rightarrow y + a$; $x \rightarrow kx$ with $k^2 = e^a$. This invariance implies that, given a solution with some value of $y(0)$, we can obtain the solution with any other value of $y(0)$ by simple rescaling. Therefore, only one of the two integration constants needed in the solution to Eq. (25.15) is really non-trivial. Hence it must be possible to reduce the degree of the equation from two to one by a judicious choice of variables. One such set of variables is:

$$v \equiv \frac{m}{x}; \quad u \equiv \frac{nx^3}{m} = \frac{nx^2}{v}. \quad (25.16)$$

In terms of v and u , Eq. (25.12) becomes

$$\frac{u}{v} \frac{dv}{du} = -\frac{(u-1)}{(u+v-3)}. \quad (25.17)$$

The boundary conditions $y(0) = y'(0) = 0$ translate into the following: v is zero at $u = 3$, and $(dv/du) = -5/3$ at $(3,0)$. (You can prove this by examining the behaviour of Eq. (25.14) near $x = 0$ retaining up to necessary order in x .)

The solution $v(u)$ to Eq. (25.17) can be easily obtained numerically: it is plotted in Fig. 25.1 as the spiraling curve. The singular points of this differential equation are given by the location in the uv plane at which both the numerator and denominator of the right hand side of Eq. (25.17) vanish together. Solving $u = 1$ and $u + v = 3$ simultaneously, we get the singular point to be $u_s = 1$, $v_s = 2$. Using Eq. (25.16), we find that this point corresponds to the asymptotic solution $n = (2/x^2)$, $m = 2x$. It is obvious from the nature of Eq. (25.17) that the solution curve will spiral around the singular point asymptotically approaching the $n = 2/x^2$ solution at large x .

The solution

The nature of the solution (shown in Fig. 25.1) allows us to put interesting bounds on various physical quantities including energy. To see this, we compute the total energy E of the isothermal sphere. The potential and kinetic energies are

*The important parameter:
 RE/GM^2*

$$\begin{aligned} U &= - \int_0^R \frac{GM(r)}{r} \frac{dM}{dr} dr = - \frac{GM_0^2}{L_0} \int_0^{x_0} mnxdx \\ K &= \frac{3}{2} \frac{M}{\beta} = \frac{3}{2} \frac{GM_0^2}{L_0} m(x_0) = \frac{GM_0^2}{L_0} \frac{3}{2} \int_0^{x_0} nx^2 dx, \end{aligned} \quad (25.18)$$

where $x_0 = R/L_0$ is the boundary and the expression for K follows from the velocity dependence of f in Eq. (25.9). The total energy is, therefore,

$$\begin{aligned} E &= K + U = \frac{GM_0^2}{2L_0} \int_0^{x_0} dx(3nx^2 - 2mnx) \\ &= \frac{GM_0^2}{2L_0} \int_0^{x_0} dx \frac{d}{dx} \{2nx^3 - 3m\} = \frac{GM_0^2}{L_0} \{n_0 x_0^3 - \frac{3}{2} m_0\}, \end{aligned} \quad (25.19)$$

where $n_0 = n(x_0)$ and $m_0 = m(x_0)$. The dimensionless quantity (RE/GM^2) is given by

$$\lambda \equiv \frac{RE}{GM^2} = \frac{1}{v_0} \{u_0 - \frac{3}{2}\}. \quad (25.20)$$

Note that the combination (RE/GM^2) is a function of only the values of (u, v) at the boundary.

A cute result

Let us now consider the constraints on λ . Suppose we specify some value for λ by specifying R, E and M . Then such an isothermal sphere

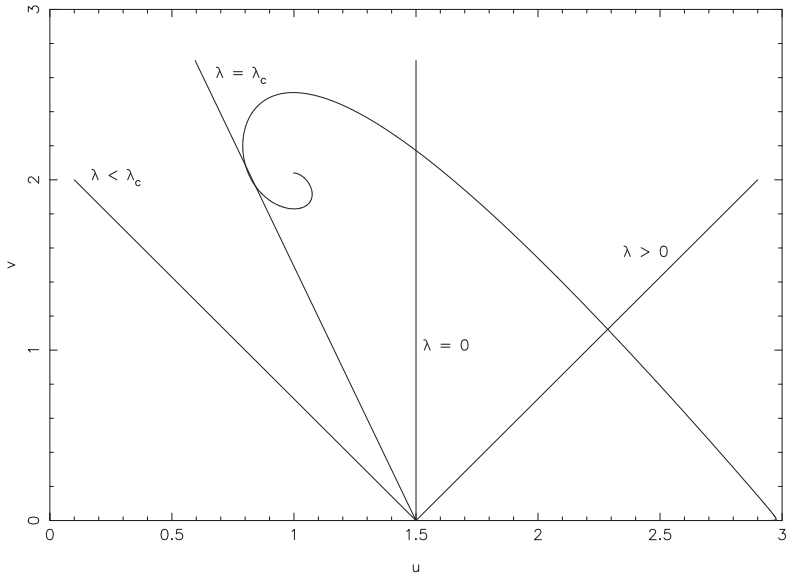


Fig. 25.1: Bound on RE/GM^2 for the isothermal sphere. See text for discussion.

must lie on the curve

$$v = \frac{1}{\lambda} \left(u - \frac{3}{2} \right); \quad \lambda \equiv \frac{RE}{GM^2}, \quad (25.21)$$

which is a straight line through the point $(1.5, 0)$ with the slope λ^{-1} . On the other hand, since *all* isothermal spheres must lie on the $u - v$ curve, *an isothermal sphere can exist only if the line in Eq. (25.21) intersects the $u - v$ curve.*

For large positive λ (positive E), there is only one intersection. When $\lambda = 0$, (zero energy) we still have a unique isothermal sphere. (For $\lambda = 0$, Eq. (25.21) represents a vertical line through $u = 3/2$.) When λ is negative (negative E), the line can cut the $u - v$ curve at more than one point; thus more than one isothermal sphere can exist with a given value of λ . (Of course, the degeneracy is lifted by specifying M, R, E individually.) But as we decrease λ (more and more negative E), the line in Eq. (25.21) will slope more and more to the left; and when λ is smaller than a critical value λ_c , the intersection will cease to exist. So we reach the key conclusion that *no isothermal sphere can exist if (RE/GM^2) is below a critical value λ_c .* This fact follows immediately from the nature of $u - v$ curve and Eq. (25.21). The value of λ_c can be found from the numerical solution and turns out to be about (-0.335) .

This result was originally due to Antonov [114] while this specific derivation was provided by me [111, 113]. It is surprising that Chan-

drasekhar, who worked out the isothermal sphere in $u - v$ coordinates as early as 1939, missed discovering the energy bound shown in Fig. 25.1. Chandrasekhar [112] has the $u - v$ curve but did not over-plot lines of constant λ . If he had done that, he would have discovered Antonov instability decades before Antonov did [114].

To understand the implications of this result, let us consider constructing such a system with a given mass M , radius R and an energy $E = -|E|$ which is negative. (The last condition means that the system is gravitationally bound.) In this case, $\lambda = RE/GM^2 = -R|E|/GM^2$ is a negative number but let us assume that it is above the critical value; that is, $\lambda > \lambda_c$. In this case we know that an isothermal sphere solution exists for the given parameter values. By construction, this solution is the local extremum of the entropy and could represent an equilibrium configuration if it is also a global maximum of entropy.

The destabilizing influence of gravity

However, for the system we are considering, it is actually quite easy to see that there is no *global* maximum for entropy. This is because, for a system of point particles interacting via the Newtonian potential, there is no lower bound to the gravitational potential energy. If we build a compact core of mass $m < M$ and radius r inside the spherical cavity, then, by decreasing r , one can supply an arbitrarily large amount of energy to the rest of the particles. Very soon, the remaining particles will have very large kinetic energy compared to their gravitational potential energy and will essentially bounce around inside the spherical cavity like a non-interacting gas of particles. The compact core in the center will continue to shrink thereby supplying energy to the rest of the particles. It is easy to see that such a core-halo configuration can have arbitrarily high values for the entropy. All this goes to show that the isothermal sphere cannot be a *global* maximum for the entropy. (This was the caveat in the calculation we performed to derive the isothermal sphere equation in Eq. (25.10) without a short distance cut-off; we tacitly assumed that the extremum condition can be satisfied for a finite value of entropy.)

No global maxima for entropy; no real equilibrium

If the radius of the spherical cavity is increased (with some fixed value for $E = -|E|$), the parameter λ will become more and more negative and for sufficiently large R , we will have a situation with $\lambda < \lambda_c$. Now the situation gets worse. The system does not even have a *local* extremum for the entropy and will evolve directly towards a core-halo configuration. This is closely related to the Antonov instability [113, 114].

Sometimes, not even a local maxima

In practice, of course, there is always a short distance cut-off because of which the core cannot shrink to an arbitrarily small radius. In such a case, there is a global maximum for entropy achieved by the (finite) core-halo configuration which could be thought of as the final state in the evolution of such a system. It will be highly inhomogeneous and, in fact, is very similar to a system which exists as a mixture of two phases. This is one key peculiarity introduced by long range attractive interactions in statistical mechanics.

Real life

The electric field lines of a point charge go out radially from it (see [figure 26.1](#)). If the charged particle is replaced by a point source of light, the photons emitted by the source *also* propagate in the same way. In other words, the electric field lines track the photon path when the light source is replaced by a charge.

Field lines from charge vs. rays of light from a source

Let us next consider a source of light kept in a gravitational field, say, near the surface of Earth. We know that light rays are bent by the action of gravity and they will no longer be propagating radially outwards. But what happens to the electric field lines of a point charge held at rest in Earth's gravitational field? Of course, we have no right to expect the simple analogy between the light source and the point charge to hold in the presence of gravity. *So it comes as a delightful surprise that it indeed holds.*

Gravity bends light rays and the field lines!

The electric field lines from a point charge — and the rays of light when the charge is replaced by a source of light — follow the same trajectory even in a constant gravitational field! They both get distorted in the same way as shown in [Fig. 26.1](#). (In fact both trajectories turn out to be arcs of circles!) This chapter is devoted to explaining this — and related — beautiful results.

Obtaining the electric field lines in the presence of Earth's gravity is a bit of a complicated task because we need to solve the Maxwell equations in a curved spacetime after first determining the form of the metric in a constant gravitational field. Given the complications, we will attack the problem in a step-by-step manner. We will first obtain the form of the relevant metric and then get the path of light rays in that metric. This will tell us how gravity bends the light rays. We will then find the electrostatic potential due to the point charge in this gravitational field (which turns out to be a rather cute result by itself). Finally, we will get the electric field lines and show that they match with the path of light rays obtained earlier.

The strategy of the chapter

A constant gravitational field, of course, is equivalent to a uniform acceleration. So the natural coordinate system for discussing a constant gravitational field \mathbf{g} is the Rindler coordinate system which can be inter-

Step 1: Metric in weak gravitational field from Principle of Equivalence

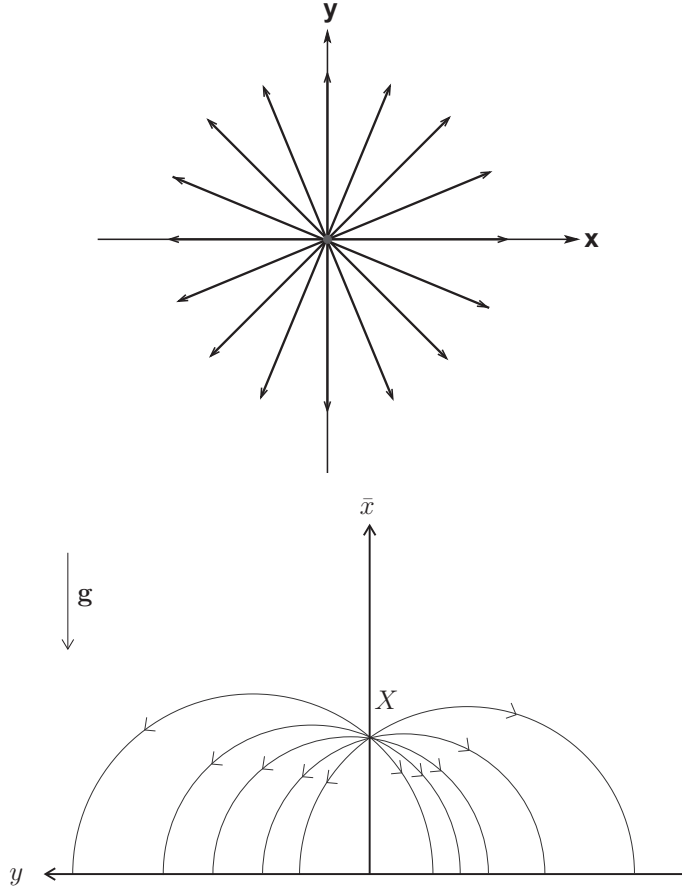


Fig. 26.1: *Top:* The field lines of a point charge in empty space extend radially outward from the charge. If we replace the charge by a point source of light, the light rays will also follow the same trajectory as the electric field lines. *Bottom:* Gravity bends the path of light rays. If a source of light is kept at the point X , the path of light rays will be as shown in the presence of a constant downward gravitational field. Incredibly enough, if the source of light is replaced by a point charge, its electric field lines will also be bent by the gravitational field in *exactly* the same manner! In other words, both the top and bottom figures can be interpreted either in terms of light rays or in terms of electric field lines.

preted in term of the coordinate system adopted by a uniformly accelerated observer in flat spacetime. The metric in the Rindler frame can be expressed in the form (see Appendix of Chapter 15):

$$\begin{aligned} ds^2 &= -(1 + \mathbf{g} \cdot \mathbf{r})^2 dt^2 + d\mathbf{r}^2 \equiv -N^2(\mathbf{r}) dt^2 + d\mathbf{r}^2 \\ &= -(1 + gx)^2 dt^2 + dx^2 + dy^2 + dz^2 . \end{aligned} \quad (26.1)$$

The second equality defines $N \equiv \sqrt{|g_{00}|}$ and the form of the metric in the third part is obtained by rotating the spatial coordinates so that the acceleration \mathbf{g} is along the x -axis. Spatial sections are flat and hence the concept of 3-vector operations in the $t = \text{constant}$ surfaces are well-defined by the usual rules of Cartesian vectors.

The transformation equations from the inertial co-ordinates (denoted by capital letters) $(T, \mathbf{R}) = (T, X, Y, Z)$, to the Rindler co-ordinates (t, x, y, z) are given by $Y = y$, $Z = z$ and

$$gT = (1 + gx) \sinh(gt); \quad 1 + gX = (1 + gx) \cosh(gt) \quad (26.2)$$

(see Appendix of Chapter 15; Eq. (15.43)). This transformation covers the quadrant $|gT| < (1 + gX)$, $(1 + gX) > 0$ of the inertial frame which will be adequate for our purpose. The transformation in Eq. (26.2) reduces to an identity (i) when $g = 0$, or (ii) at the hypersurface $t = T = 0$ even with non-zero g . On this hypersurface, $(\partial X^a / \partial x^b) = \text{dia}(N, 1, 1, 1)$. These facts are useful while transforming the tensors from one frame to another.

We will be often interested in the case of a weak acceleration and work with expressions which are accurate to first order in g . In this limit, the transformations in Eq. (26.2) reduce to

$$T \approx t(1 + \mathbf{g} \cdot \mathbf{r}); \quad \mathbf{R} \approx \mathbf{r} + (1/2)\mathbf{g}t^2. \quad (26.3)$$

The second relation is obvious; from Newtonian physics, the first one can be interpreted as the effect of gravity on the rate of clocks due to the gravitational redshift factor (see Chapter 11). These are correct to linear order in g . From Eq. (26.3), we also have the inverse transformations, again to the lowest order in g :

$$t \approx T(1 - \mathbf{g} \cdot \mathbf{R}); \quad \mathbf{r} \approx \mathbf{R} - (1/2)\mathbf{g}T^2. \quad (26.4)$$

Note that to linear order in g , we have $\mathbf{g} \cdot \mathbf{R} \simeq \mathbf{g} \cdot \mathbf{r}$. In the Rindler frame, our expressions are correct to $\mathcal{O}(\mathbf{g} \cdot \mathbf{r}/c^2)$ while in the inertial frame, they are correct to order $\mathcal{O}(\mathbf{g} \cdot \mathbf{r}/c^2)$ and $\mathcal{O}(v/c)$, where v is the speed of a particle moving with acceleration g .

When we are not interested in the $g \rightarrow 0$ limit, it is more convenient to work with a shifted x -coordinate $\bar{x} = x + g^{-1}$ in which the Rindler metric takes the form

Another coordinate system

$$ds^2 = -(g\bar{x})^2 dt^2 + d\bar{x}^2 + dy^2 + dz^2, \quad (26.5)$$

with the coordinate transformations in Eq. (26.2) becoming:

$$T = \bar{x} \sinh(gt); \quad X = \bar{x} \cosh(gt). \quad (26.6)$$

Curves of constant \bar{x} correspond to particles traveling on uniformly accelerated trajectories.

Line intervals

In this form, the transformations reduce to those corresponding to the polar coordinates if we analytically continue the time coordinates to purely imaginary values: $t \rightarrow it_E; T \rightarrow iT_E$. The proper interval between any two events in the Rindler frame can be written down just by inspection if we note that — when we use the coordinate $\bar{x} = x + g^{-1}$ and analytically continue to Euclidean space — the Euclidean distance in the plane between (t_1^E, x_1) and (t_2^E, x_2) is given by the standard cosine formula

$$s_E^2(2, 1) = x_2^2 + x_1^2 - 2x_2x_1 \cos g(t_2^E - t_1^E) . \quad (26.7)$$

Analytically continuing back and adding the transverse contribution $\rho^2 \equiv (y_2 - y_1)^2 + (z_2 - z_1)^2$, we get

$$s^2(2, 1) = \rho^2 + \bar{x}_2^2 + \bar{x}_1^2 - 2\bar{x}_2\bar{x}_1 \cosh g(t_2 - t_1) , \quad (26.8)$$

which will be useful in our discussion later on.

Step 2: Path of light rays in weak gravity

After all this background, let us determine the paths of light rays passing through any event \mathcal{P} in the Rindler frame. It can be easily verified that the paths of light rays in this xy plane are parts of circles. (Experts will note that the Rindler metric in Eq. (26.5) is conformal to a metric for which the spatial section is a Poincaré half-plane. Since the Poincaré half-plane is known to have circles as geodesics, this result is obvious. If you are not an expert, you can easily work it out!) To prove this, we begin with the generally covariant form of the Hamilton-Jacobi equation for a photon (see Eq. (2.16)), obtained by substituting $p_i = \partial_i S$ into $p_i p^i = 0$, getting:

$$g^{ik} \frac{\partial S}{\partial x^i} \frac{\partial S}{\partial x^k} = 0 , \quad (26.9)$$

where S is the action and the metric is given by Eq. (26.5). If the tangent vector to the light ray emanating from an event \mathcal{P} is $k^a = (\omega, \mathbf{k})$, then we can always choose the transverse coordinates (y, z) such that \mathbf{k} lies in the xy plane. Since we are interested in the null geodesics in the xy plane in a static metric, we can separate the variables as:

$$S = -\mathcal{E}t + yk_y + S_1(\bar{x}) , \quad (26.10)$$

where \mathcal{E} is the energy, and k_y is the y -component of the momentum and $S_1(\bar{x})$ stands for the term in the action that depends only on \bar{x} . Using Eq. (26.10) in Eq. (26.9) with Eq. (26.5) we get:

$$S = \int (\mathcal{E}^2 - k_y^2 g^2 \bar{x}^2)^{1/2} \frac{d\bar{x}}{g\bar{x}} + k_y y - \mathcal{E}t . \quad (26.11)$$

To determine the trajectory in the xy plane, we differentiate S with respect to k_y , and equate to a constant y_0 , getting

$$y - y_0 = k_y \int d\bar{x} \frac{g\bar{x}}{(\mathcal{E}^2 - k_y^2 g^2 \bar{x}^2)^{1/2}} . \quad (26.12)$$

With the substitution $k_y g \bar{x} = \mathcal{E} \cos \theta$, the above integral can be easily evaluated to find $y - y_0 = (\mathcal{E}/k_y g) \sin \theta$ so that the equation to the light ray is:

$$\bar{x}^2 + (y - y_0)^2 = R^2 , \quad (26.13)$$

where $R = \mathcal{E}/k_y g$. This is the equation to a circle (see Fig. 26.1) with center at $(\bar{x}, y) = (0, y_0)$ and having radius $R = \mathcal{E}/k_y g$.

Light rays are arcs of circles!

We also need the notion of a suitable “distance” or “time” along the light ray. This is given by a concept called the affine parameter λ which is used to parametrize light paths as $x^a(\lambda)$ just as we use proper time τ to parametrize trajectories of particles as $x^a(\tau)$. The formal definition of the affine parameter uses the machinery of general relativity which we do not want to get into. But since the metric is independent of y , we can define the affine parameter by $d^2 y / d\lambda^2 = 0$. So, with suitable initial conditions, one can simply take the y -coordinate itself as proportional to the affine parameter λ . To relate this affine parameter with the time coordinate t , we need to determine y in terms of t . Along the null trajectory, we have:

$$g^2 \bar{x}^2 dt^2 = d\bar{x}^2 + dy^2 , \quad (26.14)$$

from which we obtain

$$g^2 \bar{x}^2 \left(\frac{dt}{dy} \right)^2 = 1 + \left(\frac{d\bar{x}}{dy} \right)^2 . \quad (26.15)$$

However, from Eq. (26.13) giving the trajectory, we know that:

$$\left(\frac{d\bar{x}}{dy} \right)^2 = \left(\frac{y - y_0}{\bar{x}} \right)^2 . \quad (26.16)$$

Hence, Eq. (26.15) becomes:

$$g^2 \bar{x}^2 \left(\frac{dt}{dy} \right)^2 = 1 + \left(\frac{y - y_0}{\bar{x}} \right)^2 = \left(\frac{R}{\bar{x}} \right)^2 , \quad (26.17)$$

giving

$$t = \frac{R}{g} \int \frac{dy}{\bar{x}^2} . \quad (26.18)$$

With the substitutions $y = y_0 + R \sin \theta$ and $\bar{x} = R \cos \theta$, the above integral can be evaluated to give:

$$t = \frac{1}{g} \log \left(\frac{1 + \tan(\theta/2)}{1 - \tan(\theta/2)} \right) = \frac{1}{g} \log \tan \left(\frac{\theta}{2} + \frac{\pi}{4} \right). \quad (26.19)$$

Substituting back in terms of the original variables, we find:

$$2 \tanh^{-1}(e^{gt}) - \frac{\pi}{2} = \sin^{-1} \left(\frac{y - y_0}{R} \right). \quad (26.20)$$

Rearranging and simplifying, we have:

$$\tanh gt = \left(\frac{y - y_0}{R} \right) = \frac{\alpha}{R} \lambda. \quad (26.21)$$

Therefore, the affine parameter turns out to be proportional to $\tanh gt$, where α is a proportionality constant giving $y = y_0 + \alpha \lambda$. We fix it by noting that when $g \rightarrow 0$, we would like the affine parameter to become t . This gives $\lambda = g^{-1} \tanh(gt)$.

Step 3: Electrostatic potential in weak gravity

We are now in a position to determine the electrostatic potential and the electric field of a charged particle which is at rest at the origin of the Rindler frame or — equivalently — in a weak homogeneous gravitational field. Such a charged particle will be moving along a uniformly accelerated trajectory in the inertial coordinate system.

We begin by noting that, because of the static nature of the Rindler frame, the four-vector potential reduces to the form $A_i = (A_0, 0, 0, 0)$ with $A_0(\mathbf{r})$ being independent of the time coordinate. It is therefore enough if we determine the electrostatic potential on the $t = 0$ hypersurface. We also know that the potential at an event x^i is determined by the nature of the trajectory of the charged particle $z^i(t_R)$ at the retarded time t_R . This retarded time is a function of the field coordinates x^i and is determined by the condition that $z^i(t_R)$ and x^i are connected by a light ray. We will argue that the potential $A_0(0, \mathbf{r})$ due to a charge at rest in the Rindler frame should be expressible in the form

$$A_0(\mathbf{r}) = A_0(0, \mathbf{r}) = \frac{q}{\lambda(\mathcal{F}; \mathcal{S})}, \quad (26.22)$$

where $\lambda(\mathcal{F}; \mathcal{S})$ is the affine parameter distance along a null geodesic connecting the field event $\mathcal{F}(0, \mathbf{r})$ with the location of the source at the retarded time $\mathcal{S}(t_R, \mathbf{0})$. So it is *just like a Coulomb field with the affine parameter distance replacing the radial distance*.

The proof

This result is easily established along the following lines: We begin with the usual formula for the potential of an arbitrarily moving charge in

It is just Coulomb potential with affine parameter for distance!

inertial coordinates, written in the form (see Appendix):

$$A_k = \frac{2qu_k}{|ds^2/d\tau|}, \quad (26.23)$$

where $u^i(\tau)$ is the four-velocity of the charge in the inertial frame at the proper time τ and the expression on the right hand side has to be evaluated at the retarded time on the trajectory of the charge. Taking the dot product of both sides with u^k (at the retarded time) we get the scalar equation $A_k u^k = -2q/|ds^2/d\tau|$. In the Lorentz frame in which the charge was at rest at the origin, at the retarded time, the right hand side reduces to usual Coulomb form $q/|T_R|$ where T_R (< 0) is the relevant retarded time satisfying the condition $T_R = -|\mathbf{R}|$. We next note that $-T_R$ or $|\mathbf{R}|$ is actually the affine distance λ along the null geodesic connecting the event $\mathcal{S} = (T_R, \mathbf{0})$ corresponding to the source at retarded time to the event $\mathcal{F} = (0, \mathbf{R})$ where the field is measured. This shows that we can equivalently write $A_k u^k = -q/\lambda$ in any Lorentz frame, for an arbitrarily moving charged particle. But both sides of this equation are also generally covariant in flat spacetime when curvilinear coordinates are used. (As an aside, let me make the following comment: In the left hand side, A_k is the potential at some event x^i while u^k is the four-velocity of the charge at the retarded event z^i connected to x^i by a light ray. So the dot product of these two vectors, defined *a priori* in two different events, can be taken only after parallel transporting one vector to the location of another. Since this parallel transport is unique in flat spacetime, the expression is invariant with respect to curvilinear coordinate transformations in flat spacetime. Unfortunately, this prevents us from applying this idea to genuinely curved spacetime without modification.) Therefore we can use the same relation in curvilinear coordinates as well, and express the electrostatic potential of a static source at the origin of the Rindler frame in a generally covariant manner, in terms of the affine parameter distance between the source at the retarded time and the field point. In the Rindler frame, we have $A_k u^k = A_0$ since $u^0 = 1/N = 1$ at the trajectory of the charge at all times, including at the relevant retarded time, thereby leading to the result in Eq. (26.22). Since the affine parameter is given by Eq. (26.21), we get the result:

$$A_0(\mathbf{r}) = \frac{q}{g^{-1} \tanh g t_R}. \quad (26.24)$$

Obviously, both λ and the retarded time t_R depend on the spatial coordinate of the field point \mathbf{r} . So we need to next compute the retarded time t_R . Consider the field event $\mathcal{F} = (0, \mathbf{r})$ and the source event at the retarded time $\mathcal{S} = (t_R, \mathbf{0})$, connected by a null ray. Setting $s^2 = 0$ in the expression for the interval given by Eq. (26.8) will allow us to determine t_R . In

Why we can't use this in general

Eq. (26.8), we are now interested in the case with $\bar{x}_1 = g^{-1}$, $y_1 = z_1 = 0$, $t_1 = t_R$, $t_2 = 0$, $\bar{\mathbf{r}}_2 = \bar{\mathbf{r}}$ for which we get:

$$s^2(\mathcal{F}; \mathcal{S}) = \rho^2 + \bar{x}^2 + g^{-2} - 2g^{-1}\bar{x} \cosh gt_R. \quad (26.25)$$

The condition $s^2 = 0$ now determines t_R in terms of other variables and we get:

$$\cosh gt_R = \frac{g}{2\bar{x}} [\rho^2 + \bar{x}^2 + g^{-2}]. \quad (26.26)$$

More explicitly, we have

$$\cosh gt_R = \frac{1 + g^2 \bar{r}^2}{2g\bar{x}}; \quad \sinh gt_R = \frac{1}{2g\bar{x}} [(1 + g^2 \bar{r}^2)^2 - 4g^2 \bar{x}^2]^{1/2}, \quad (26.27)$$

where $r^2 = x^2 + y^2 + z^2 \equiv \rho^2 + x^2$. Taking the ratio to obtain $\tanh gt_R$ and switching back to $x = \bar{x} - g^{-1}$, leads to the expression for the potential given by

Final result

$$A_0 = \frac{q}{r} \frac{1 + gx + g^2 r^2 / 2}{(1 + gx + g^2 r^2 / 4)^{1/2}}. \quad (26.28)$$

While this expression has been obtained by several people in the past [115–121], the cute interpretation in terms of the affine parameter in Eq. (26.22) is from Ref. [122]. This result can also be expressed [123] as

$$A_0 = \frac{qg}{2} \left(\frac{\ell_+}{\ell_-} + \frac{\ell_-}{\ell_+} \right); \quad \ell_{\pm}^2 = \rho^2 + (\bar{x} \pm g^{-1})^2, \quad (26.29)$$

where ℓ_{\pm} represent the distances to the field point from a charge (at $1/g$) and an ‘image charge’ (at $-1/g$). Equipotential surfaces correspond to constant values of ℓ_+/ℓ_- . Since the locus of a point that moves keeping the ratio of distances from two different points constant, is a circle, we find that equipotential surfaces are circles in the xy plane.

To get the electric field from the vector potential, we note that when the charge distribution is static we can assume that only A_0 and $F_{\mu 0} = -F_{0\mu} = \partial_{\mu} A_0$ are non-zero, leading to

$$\mathbf{E} = -N^{-1} \nabla A_0. \quad (26.30)$$

So, in the Rindler frame with the electric field given by:

$$\mathbf{E} = -\frac{\nabla A_0}{(1 + \mathbf{g} \cdot \mathbf{r})}. \quad (26.31)$$

Without loss of generality, we can confine our attention to the xy plane with $\mathbf{E} = (E^x, E^y, 0)$. Explicit calculation using Eq. (26.28) gives:

$$E_x = \frac{qx}{r^3} \frac{1 + gx/2 - gy^2/2x}{(1 + gx + g^2 r^2/4)^{3/2}}; \quad E_y = \frac{qy}{r^3} \frac{1 + gx}{(1 + gx + g^2 r^2/4)^{3/2}}. \quad (26.32)$$

Since only A_0 is non-zero in the Rindler frame, it follows trivially that the magnetic field vanishes.

We can now obtain our final result, related to the bending of the electric field lines by gravity, directly from this expression. We know that the electric field lines in the xy plane are given by curves $x = x(y)$ which satisfy the equation $dx/dy = E^x/E^y$. On using Eq. (26.32), this reduces to

$$\frac{dx}{dy} = \frac{(x + g^{-1})^2 - g^{-2} - y^2}{2y(x + g^{-1})}. \quad (26.33)$$

It is easy to verify that this equation is solved by the circles in Eq. (26.13) by noting that, for these circles, Eq. (26.33) gives $dx/dy = -(y - y_c)/(x + g^{-1})$ which is the same relation we get from Eq. (26.13). In other words, *the electric field lines of a static charge in the Rindler frame coincide with the paths of light rays!* It is understandable that the electric field lines bend under the action of gravity but it is rather surprising that they do so exactly like the light rays (Fig. 26.1).

Step 4: The electric field lines in weak gravity

Field lines behave like light rays!

Having obtained the exact results, we shall next consider the case of a *weak* gravitational field and work out the expressions to the linear order in g . (A Rindler frame with acceleration \mathbf{g} corresponds to a weak gravitational field $-\mathbf{g}$ in the direction *opposite* to the acceleration; but for simplicity, we shall continue to quote the results in terms of \mathbf{g} .) In this case, we get the solution:

$$A_0 = \frac{q}{r} \left(1 + \frac{\mathbf{g} \cdot \mathbf{r}}{2} \right). \quad (26.34)$$

This is the electrostatic potential, in the limit of a weak gravitational field, of a charge at rest at the origin of co-ordinates in the Rindler frame. We can use Eq. (26.31) to obtain the corresponding electric field from this potential. We get:

$$\mathbf{E} = \frac{q\hat{\mathbf{r}}}{r^2} - \frac{q}{2r} (\mathbf{g} + (\mathbf{g} \cdot \hat{\mathbf{r}})\hat{\mathbf{r}}) = \frac{q\hat{\mathbf{r}}}{r^2} \left(1 - \frac{(\mathbf{g} \cdot \mathbf{r})}{2} \right) + \frac{q}{2r} (-\mathbf{g}), \quad (26.35)$$

where $\hat{\mathbf{r}}$ denotes the unit vector in the radial direction. In the first expression for \mathbf{E} in Eq. (26.35), we have given the result in terms of a Coulomb term plus a correction due to the gravitational field. In the second expression, we have separated the two terms based on the direction of the vectors: the first one is in the radial direction with a corrected Coulomb term while the second one is in the direction of the gravitational field ($-\mathbf{g}$).

These results are for a charge located at the origin of the Rindler frame. For our next application, we will require the potential and field produced by a charge at rest, not at the origin, but at an arbitrary point $\mathbf{r}_0 = (x_0, y_0, 0)$. (As noted before, there is no loss of generality in confining to the xy plane.) It is not obvious that we can simply introduce a translation of coordinates because our background metric is not translationally invariant. What *is* surprising, however, is that the electric field

does turn out to be translationally invariant to linear order in g . (For a rigorous proof, see [122].) We find:

$$\mathbf{E} = \frac{q\boldsymbol{\ell}}{\ell^3} - \frac{q}{2\ell}(\mathbf{g} + (\mathbf{g} \cdot \hat{\boldsymbol{\ell}})\hat{\boldsymbol{\ell}}); \qquad \boldsymbol{\ell} = \mathbf{r} - \mathbf{r}_0, \qquad (26.36)$$

so that this electric field depends only on the vectorial separation between the charge and the field point.

The results obtained above lead to an interesting consequence when we consider the forces exerted by two charges — located in a weak gravitational field — on each other. To provide a concrete realization of this situation, consider the following thought experiment. Two charged particles of masses m_1 and m_2 and charges q_1 and q_2 are held supported in a weak gravitational field by, for example, hanging the two particles by strings attached to the ceiling of a room in Earth’s gravitational field, so that the charges are located on the same horizontal plane (Fig. 26.2). If the particles were uncharged, the sum of the tensions on the two strings will be

*The weight of
electrostatic energy*

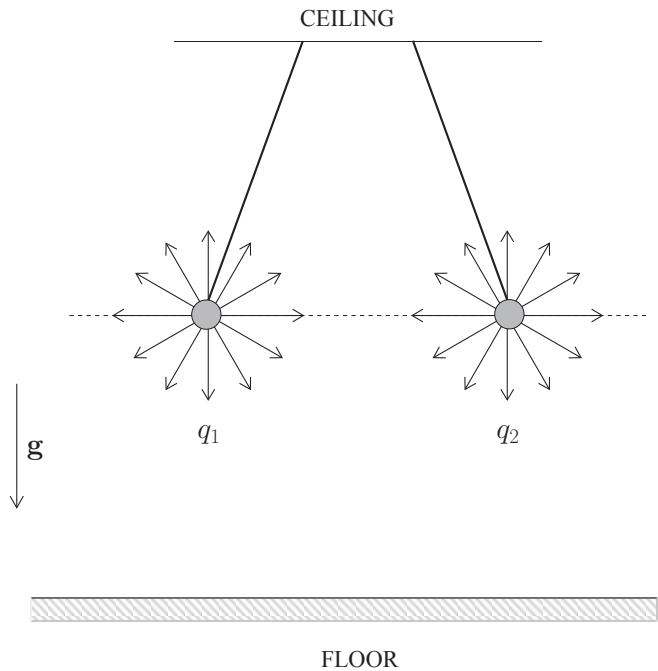


Fig. 26.2: Two charged particles are held supported in a weak gravitational field by hanging them by strings attached to the ceiling of a room in Earth’s gravitational field. *If we ignore the effect of gravity on the electrostatic field produced by the charges, the force exerted by the charges on one another is the usual Coulomb force, directed horizontally along the line joining the charges. They cancel each other and there is no net electrostatic force acting on the charges.*

equal to the total weight of the particles, $(m_1 + m_2)g$. When the particles are charged, they exert electrostatic forces on one another. If we ignore the effect of gravity on the electrostatic field produced by the charges, then the force exerted by the charges on one another is the usual Coulomb force which is directed horizontally along the line joining the charges. These Coulomb forces cancel each other and there is no *net* electrostatic force acting on the charges.

The situation changes in a curious manner when we take into account the distortion of the field lines due to the weak gravitational field (Fig. 26.2). From Eq. (26.36) we find that there is a component of the electric field in the direction of $-\mathbf{g}$ produced by each charge at the location of the other. When we add up the forces exerted by the two charges on each other, the forces in the direction of ℓ cancel out leading to the net extra force given by

$$\mathbf{F}_{12} + \mathbf{F}_{21} = -\frac{q_1 q_2}{\ell} \mathbf{g} = \frac{q_1 q_2}{\ell} \mathbf{g}_e. \quad (26.37)$$

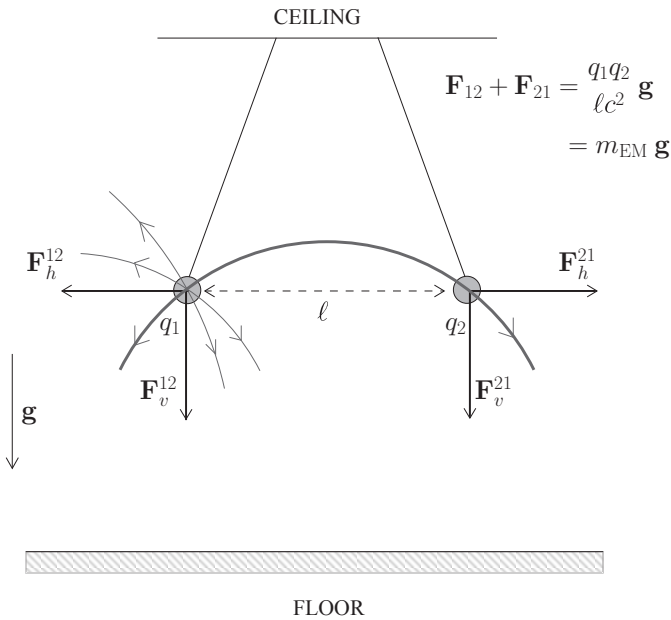


Fig. 26.3: The same situation as shown in the previous figure. Now we take into account the distortion of the field lines due to the weak gravitational field, which produces a component of the electric field in the downward direction of at the location of each charge. While the horizontal forces cancel, these vertical components add up. The two strings supporting the charges now have to support an additional weight $(q_1 q_2 / \ell c^2)g$ which is the weight of the electrostatic potential energy in this frame. In a freely falling frame these charges are moving with an acceleration g and this force will be interpreted as due to the radiation field.

In the last expression we have used the fact that the direction of acceleration in the Rindler frame \mathbf{g} and the direction of Earth's gravitational field \mathbf{g}_e are opposite to one another. This result shows that the two strings supporting the charges located in a weak gravitational field have to support an additional weight $(q_1 q_2 / \ell c^2)g$ which can be interpreted as the weight of the electrostatic potential energy. In fact, we can turn this argument around to claim that the distortion of the electric field due to gravity *must* produce a term of the form $(q/\ell)\mathbf{g}$ since gravity has to support the electrostatic energy. Obviously, the result extends to any number of charged particles all located in the same horizontal plane; the extra weight that needs to be supported by the string will be equal to the effective weight of the total electrostatic energy of the system. (It appears that this problem was first tackled by Enrico Fermi in ref. [124]. Subsequently, there have been several papers exploring this issue the results of which did not always agree with each other; see e.g., Refs. [125, 126]. These papers also contain more extensive bibliography.)

*A non-trivial
application:
radiation reaction
on the charge*

Finally, we shall consider an intriguing application of the above analysis: that of determining the *radiation reaction force* on an accelerated charged particle. We know that a charge with *variable* acceleration will feel a radiation reaction force in the inertial frame proportional to \dot{g} in the non-relativistic limit. In Chapter 20 we argued (based on [127]) that the electromagnetic fields of this charge — with *variable* acceleration — can actually be determined from knowing only the fields of a uniformly accelerated charge. The question arises as to whether we can also interpret the radiation reaction in the Rindler frame. We will now derive this result in the non-relativistic limit.

We know that a charged particle which has a uniform acceleration \mathbf{g} in the inertial frame can be mapped to a charged particle at rest at the origin of the Rindler frame. The electric field produced by this charge in the Rindler frame is given by Eq. (26.35) which is accurate to lowest order in g . Since the Rindler frame is a static frame of reference, we can, without loss of generality, choose to measure this field at the time $t = 0$.

Let us now suppose that the charged particle is at the origin of the inertial frame (which coincides with the origin of the Rindler frame) at $t = T = 0$, but its acceleration g is slowly varying in time with a small but non-zero time derivative, \dot{g} . In other words, the instantaneous acceleration of the charged particle at any time t (near $t = 0$) can be expressed as $g(t) \approx g_0 + \dot{g}t$ where g_0 is a constant and \dot{g} is small and higher derivatives (\ddot{g} , \dddot{g} etc.) are ignored. The trajectory of this charged particle can now be expressed in the Rindler frame. The charge is no longer stationary at the origin, but has a trajectory given by $x_0(t) = \dot{g}t^3/6$. So, the position of the particle in the Rindler frame now changes with time due to the time derivative of the acceleration \dot{g} .

The self-force

It is precisely this case that we are interested in for the radiation reaction calculation. We will now derive the expression for the electric field of

a charged particle which moves with slowly varying g as described above, retaining only terms to lowest order in g throughout the analysis: First, we will obtain the expression for the electric field, along the x -axis, of a charge that is at rest at the origin of the Rindler frame. Then, we will modify this expression for the case of a charge that is not exactly at rest, but has a small but non-zero \dot{g} . To the lowest order of approximation, this can be accomplished by replacing g everywhere in the electric field expression, by $g(t) = g_0 + \dot{g}t$, and at the same time replacing x by $x - \dot{g}t^3/6$. This latter replacement is necessary because our electric field expression gives the field at point x , produced by a charge located at the origin. Since the charge now has the trajectory $x_0(t) = \dot{g}t^3/6$, the translational invariance of the field requires the replacement of x wherever it appears in the electric field expression, by $x - x_0 = x - \dot{g}t^3/6$.

We will now carry out the above procedure. Consider the electric field in the Rindler frame of a charged particle which is at rest at the origin of this frame along the x -axis of the Rindler frame. This electric field is given by setting $r = x$ in the general expression in Eq. (26.35) leading to:

$$E_x = \frac{q}{x^2} - \frac{q\dot{g}}{x}; \quad E_y = 0. \quad (26.38)$$

Replacing g by $g_0 + \dot{g}t$ and x by $x - \dot{g}t^3/6$, we get the field due to a charge with a slowly varying acceleration:

$$E_x = \frac{q}{(x - \dot{g}t^3/6)^2} - \frac{q(g_0 + \dot{g}t)}{(x - \dot{g}t^3/6)}. \quad (26.39)$$

This expression is, in general, time-dependent and has to be evaluated at the retarded time corresponding to the field point x . Again, to the lowest order of approximation, the exact nature of the curved path of the light ray does not matter and it can be approximated by a straight line connecting the point (t, x) with approximately the origin (since \dot{g} is small). Hence, we have $x^2 = t^2$, since the path of light is a null line connecting the above two points. However, since we are measuring the fields at the point $x > 0$, say, at the time $t = 0$, the retarded time is negative with $t = -x$. Effecting this substitution in Eq. (26.39) and retaining terms to lowest order in g , we obtain (what will turn out to be) a miraculous result:

$$E_x = \frac{q}{x^2} - \frac{qg_0}{x} + \frac{2}{3}q\dot{g}. \quad (26.40)$$

This expression, in the limit of $x \rightarrow 0$ is *identical to the expression for the self-force on a charge obtained by Dirac [128, 129] with exactly the same coefficients, relative signs and the nature of divergent terms!*

The first two terms are well-known divergences when $x \rightarrow 0$, (and are discussed extensively in the literature). Briefly, the first term is discarded as the electrostatic self energy and the second term, when moved to the

*Dirac's result,
derived simply!*

left hand side of the equations of motion, leads to a mass renormalization because it is proportional to the acceleration. It is interesting that, even with all our approximations — working things out to only the lowest order in g , and neglecting all higher powers of g throughout the analysis — we obtain these two terms with their appropriate signs and the correct numerical coefficients in front.

The real strength of our simple technique, however, is brought out by the production of the last term which is *identical to the standard expression for the radiation reaction field* of a charged particle. (The radiation reaction *force* will be q times this *field*, $(2/3)q^2\dot{g}$.) Again, the factor and sign in this term are identical to those in the standard expression. This computation of the radiation reaction illustrates the power of our simple non-relativistic approximation to the electric field.

Appendix: The Eq. (26.23) can be obtained as follows. We start with two standard results: (i) In the Lorentz gauge, the vector potential — related to the current by $\square A_m = -4\pi j_m$, has the solution:

$$A_m(x) = 4\pi \int d^4x' G_{ret}(x-x') j_m(x') , \quad (26.41)$$

where G_{ret} is the retarded Green's function given by:

$$G_{ret}[x] = \frac{1}{2\pi} \delta(s^2) \theta(x^0); \quad s^2 \equiv x^\mu x_\mu . \quad (26.42)$$

The $\delta(s^2)$ factor is obvious from the propagation along the lightcone since this is the only functional form which will lead to the correct $1/r$ dependence in the static case; the $\theta(x^0)$ ensures that the retarded condition is satisfied. The proportionality constant can be determined by considering the Poisson equation in the static limit. (ii) The current $j^m(x)$ for a point charge moving along a worldline $z^m(\tau)$ with a 4-velocity $u^m(\tau)$ is given by:

$$j^m(x) = e \int d\tau \delta[x - z(\tau)] u^m(\tau) . \quad (26.43)$$

This makes the current density zero everywhere except on the worldline and, on the worldline it reduces to the standard expression, if we convert the τ integration to a t integration. Equations (26.43) and (26.42) together give the vector potential to be:

$$\begin{aligned} A_m(x) &= 4\pi e \int d\tau G_{ret}[x - z(\tau)] u_m(\tau) \\ &= 2e \int d\tau \delta(s^2) u_m(\tau) = \frac{2e u_m}{ds^2/d\tau} , \end{aligned} \quad (26.44)$$

evaluated at the retarded time. This is the expression used in the main text.

References

1. K. Kuchar (1980), *Gravitation, geometry, and nonrelativistic quantum theory*, Phys. Rev. **D 22**, 1285.
2. H. Padmanabhan and T. Padmanabhan (2011), *Nonrelativistic limit of quantum field theory in inertial and non-inertial frames and the principle of equivalence*, Phys. Rev. **D 84**, 085018.
3. W. B. Case (2008), *Wigner functions and Weyl transforms for pedestrians*, Am. J. Phys., **76**, 937.
4. J. Hermann (1710), *Extrait d'une lettre de M. Herman a M. Bernoulli datee de Padoue le 12.Juillet 1710*, Histoire de l'academie royale des sciences (Paris), **1732**, 519.
5. J. Bernoulli (1710), *Extrait de la Reponse de M. Bernoulli a M. Herman datee de Basle le 7. Octobre 1710*, Histoire de l'academie royale des sciences (Paris), **1732**, 521.
6. P.S. Laplace (1799), *Traite de mecanique celeste*. Tome I Premiere Partie, Livre II, pp. 165ff.
7. W.R. Hamilton (1847), *Applications of Quaternions to Some Dynamical Questions*, Proceedings of the Royal Irish Academy **3**, page xxxvi: (Appendix III).
8. J.W. Gibbs and E.B. Wilson (1901), *Vector Analysis*, (Yale University Press, US), p. 135.
9. C. Runge (1919), *Vecktoranalysis*, (Hirzel, Leipzig) Volume I.
10. W. Lenz (1924) *Über den Bewegungsverlauf und Quantenzustände der gestörten Keplerbewegung*, Zeitschrift für Physik, **24**, 197.
11. P. G. Tait and W. J. Steele (1900), *A Treatise on the Dynamics of a Particle*, (reprinted by Adamant Media Corporation in 2005).
12. T. Padmanabhan (2009), *Perturbing Coulomb to Avoid Accidents!* Resonance, **14**, 622.
13. W. Greiner and B. Muller (1989), *Quantum mechanics — Symmetries*, Chapter 14, (Springer, Berlin).
14. F.H.J. Cornish (1984), *The hydrogen atom and the four-dimensional harmonic oscillator*, J. Phys. A, **17**, 323.
15. R. de Lima Rodrigues (2009), *On the Hydrogen Atom via Wigner-Heisenberg Algebra*, J.Phys.A, **42**, 355213.
16. C. Kacser (1959), *Higher Born approximations in non-relativistic Coulomb scattering*, Il Nuovo Cimento, **XIII**, 303.
17. W. Thompson and P.G. Tait (1962), *Principles of Mechanics and Dynamics*, (Dover, New York).
18. W.D. Macmillan (1958), *Theory of the Potential*, (Dover, New York).

19. T. Padmanabhan (2008), *Potentials of potatoes: A surprise in Newtonian Gravity*, Resonance, **13**, 4.
20. T. Padmanabhan (1996), *Cosmology and Astrophysics through Problems*, (Cambridge University Press, Cambridge).
21. T. Padmanabhan (2009), *Lagrange has (more than) a Point!* Resonance, **14**, 318.
22. R. Grenberg and D. R. Davis (1978), *Stability at potential maxima: The L_4 and L_5 points of the restricted three body problem*, Am. J. Phys., **46**, 1068.
23. D. Boccaletti and G. Pucacco (1996), *Theory of Orbits, Vol I*, page 271, (Springer, Berlin).
24. T. Padmanabhan (2009), *Extreme Physics*, Resonance, **14**, 907.
25. C. E. Mungan and T. C. Lipscombe (2013), *Complementary curves of descent*, Eur. J. Phys., **34**, 59.
26. V. Perlick (1991), *The brachistochrone problem in a stationary space-time*, J. Math. Phys., **32**, 3148.
27. G.J. Tee (1998), *Isochrones and Brachistochrones*, (Department of Mathematics, University of Auckland).
28. C. Boyer (1987), *The rainbow: From Myth to Mathematics*, (Princeton University Press, Princeton).
29. R. Goldstein (1994), *Strange Attractors: Stories*, (Penguin Books, USA).
30. M. Grossmann, E. Schmidt and A. Haussmann (2011), *Photographic evidence for the third-order rainbow*, Applied Optics, **50**, p. F134.
31. T. Padmanabhan (2008), *Ambiguities in Fluid Flow*, Resonance, **13**, 802.
32. L.D.Landau, E.M. Lifshitz (1987), *Fluid Mechanics*, Sections 10, 11 (Pergamon Press, Oxford).
33. T. E. Faber (1995), *Fluid mechanics for Physicists*, Section 4.8 (Cambridge University Press, Cambridge).
34. T. Padmanabhan (2008), *Isochronous Potentials*, Resonance **13**, 998.
35. A. B. Pippard (1989), *The Physics of Vibration*, page 15 (Omnibus Edition, Cambridge University Press, Cambridge).
36. L. D. Landau and E. M. Lifshitz (1976), *Mechanics*, (Pergammon Press, Oxford).
37. M. Asorey, J.F. Carinena, G. Marmo and A. Perelomov (2007), *Isoperiodic classical systems and their quantum counterparts*, Annals Phys., **322**, 1444.
38. R.E. Langer (1934), *The asymptotic solutions of ordinary linear differential equations of the second order, with special reference to the Stokes phenomenon*, Bull. Am. Math. Soc., **40**, 545.
39. R.E. Langer (1937), *On the Connection Formulas and the Solutions of the Wave Equation*, Phys. Rev., **51**, 669.
40. M.V. Berry and K.E. Mount (1972), *Semiclassical approximations in wave mechanics*, Rep. Prog. Phys., **35**, 315.
41. T. Padmanabhan (2008), *The Logarithms of Physics*, Resonance, **13**, 510.
42. L. R. Mead and J. Godines (1991), *An analytic example of renormalization in two dimensional quantum mechanics*, Am. J. Phys., **59**, 935.
43. P. Gosdzinsky and R. Tarrach (1991), *Learning quantum field theory from elementary quantum mechanics*, Am. J. Phys., **59**, 70.
44. B.R. Holstein (1993), *Anomalies for pedestrians*, Am. J. Phys., **61**, 142.
45. A. Cabo, J.L. Lucio and H. Mercado (1998), *On scale invariance and anomalies in quantum mechanics*, Am. J. Phys., **66**, 240.
46. M. Hans (1983), *An electrostatic example to illustrate dimensional regularization and renormalization group technique*, Am. J. Phys., **51**, 694.
47. S. A. Coon and B. R. Holstein (2002), *Anomalies in quantum mechanics: The $1/r^2$ potential*, Am. J. Phys., **70**, 513.
48. M. Visser (2005), *Heuristic approach to the Schwarzschild geometry*, Int.J.Mod.Phys. **D14**, 2051 [gr-qc/0309072].
49. T. Padmanabhan (2008), *Schwarzschild Metric at a discounted price*, Resonance, **13**, 312.

50. C.W. Misner, K.S. Thorne and J.A. Wheeler (1973), *Gravitation*, Chapter 41, (Freeman, New York).
51. T. Padmanabhan (2006), *An Invitation to Astrophysics*, Chapter 5, (World Scientific, Singapore).
52. T. Padmanabhan (2008), *Why are black holes hot?*, Resonance, **13**, 412.
53. J.A. Wheeler (1999), *A Journey into Gravity and Spacetime*, Scientific American Library, (W. H. Freeman, New York).
54. K. S. Thorne (1995), *Black Holes and Time Warps: Einstein's Outrageous Legacy*, (W. W. Norton, New York).
55. T. Padmanabhan (2002), *Classical and quantum thermodynamics of horizons in spherically symmetric spacetimes*, Class.Quan.Grav., **19**, 5387 [gr-qc/0204019].
56. T. Padmanabhan (2010), *Thermodynamical Aspects of Gravity: New insights*, Reports in Progress of Physics, **73**, 046901 [arXiv:0911.5004].
57. T. Padmanabhan (2008), *Thomas Precession*, Resonance, **13**, 610.
58. R. A. Muller (1992), *Thomas precession: Where is the torque?* Am.J.Phys., **60**, 313.
59. T. Padmanabhan (2000) *Theoretical Astrophysics, Volume 1: Astrophysical Processes*, (Cambridge University Press).
60. T. Padmanabhan (2008), *Foucault meets Thomas*, Resonance, **13**, 706.
61. J. B. Hart et al. (1987), *A simple geometric model for visualizing the motion of a Foucault pendulum*, Am. J. Phys., **55**, 67.
62. J. von Bergmann, H. von Bergmann (2007), *Foucault pendulum through basic geometry*, Am. J. Phys., **75**, 888.
63. M.I. Krivoruchenko (2009), *Rotation of the swing plane of Foucault's pendulum and Thomas spin precession: Two faces of one coin*, Phys.Usp., **52**, 821. [arXiv:0805.1136v1].
64. L. D. Landau and E. M. Lifshitz (1977), *Quantum mechanics*, p. 76, [Pergamon Press, Oxford; Third Edition].
65. T. Padmanabhan (2010), *Gravitation: Foundations and Frontiers*, p. 126, (Cambridge University Press, Cambridge).
66. T. Padmanabhan (2008), *Paraxial Optics and Lenses*, Resonance, **13**, 1098.
67. T. Padmanabhan (2009), *The Optics of Particles*, Resonance, **14**, 8.
68. T. Padmanabhan (1994), *Path integral for the relativistic particle and harmonic oscillators*, Found.Physics, **24**, 1543.
69. T. Padmanabhan (2009), *Real Effects from Imaginary Time*, Resonance, **14**, 1060.
70. K. Srinivasan and T. Padmanabhan (1999), *Particle Production and Complex Path Analysis*, Phys. Rev., **D 60**, 24007.
71. B. R. Holstein (1984), *Semiclassical treatment of above barrier scattering*, Am. J. Phys., **52**, 321.
72. T. Padmanabhan (2009), *The Power of Nothing*, Resonance, **14**, 179.
73. H.B.G. Casimir (1948), Proc. Kon. Ned. Akad. Wetensch. **B51**, 793.
74. S. K. Lamoreaux (1997), *Demonstration of the Casimir Force in the 0.6 to 6 μm Range*, Phys. Rev. Lett. **78**, 5.
75. G. Bressi, G. Carugno, R. Onofrio and G. Ruoso (2002), *Measurement of the Casimir Force between Parallel Metallic Surfaces*, Phys. Rev. Lett. **88**, 041804.
76. T. Padmanabhan (2009), *Why does an Accelerate Charge Radiate?* Resonance, **14**, 499.
77. J.J. Thomson (1907), *Electricity and Matter*, Chapter III, (Archibald Constable, London).
78. E.M. Purcell (2008), *Electricity and Magnetism*, The Berkeley Physics Course, Volume 2, 2nd ed. (Mc-Graw-Hill, New York).
79. F. S. Crawford (2008), *Waves*, The Berkeley Physics Course, Volume 3, 2nd ed. (Mc-Graw-Hill, New York).
80. H. Padmanabhan (2009), *A Simple derivation of the electromagnetic field of an arbitrarily moving charge*, Am.J.Phys., **77**, 151 [arXiv:0810.4246].

81. J. J. Thomson (1903), *The Magnetic Properties of Systems of Corpuscles describing Circular Orbits*, Phil. Mag, **45**, 673.
82. G. A. Schott (1912), *Electromagnetic Radiation* (Cambridge University Press, Cambridge).
83. L. Arzimovitch and I. Pomeranchuk (1945), *The Radiation of Fast Electrons in the Magnetic Field*, J. Phys. (USSR) **9**, 267.
84. J. Schwinger (1949), *On the Classical Radiation of Accelerated Electrons*, Phys. Rev. **75**, 1912.
85. R. P. Feynman, R.B. Leighton and M. Sands (1964) *Feynman Lectures in Physics*, Volume II; section 17-4 (Addison Wesley, USA).
86. J. M. Aguirregabiria and A. Hernandez (1981), *The Feynman paradox revisited*, Eur.J.Phys., **2**, 168.
87. T. Padmanabhan (2008), *Angular momentum of electromagnetic field*, Resonance, **13**, 108.
88. B. Hughes (1995), *Random Walks and Random Environments*, Vol I, (Oxford university press, Oxford).
89. E. W. Montroll and M.F. Shlesinger (1984), *On the wonderful world of random walks*, in *Studies in Statistical Mechanics*, edited by J.L., Lebowitz and E.W. Montroll , Vol. 11, (North-Holland, Amsterdam).
90. J. Rudnick and G. Gaspari (2004), *Elements of the Random Walk*, (Cambridge University Press, Cambridge).
91. T. Padmanabhan (2009), *Random Walk Through Random Walks - I*, Resonance, **14**, 638.
92. T. Padmanabhan (2011), *Statistical mechanics of gravitating systems and some curious history of Chandras rare misses!*, Pramana, **77** 147156.
93. S. Chandrasekhar (1942), *Principles of Stellar Dynamics*, (Dover, New York).
94. J. H. Jeans (1929), *Astronomy and Cosmogony*, (Cambridge University Press, Cambridge).
95. V. A. Ambartsumian, Uch. Zap. L.G.V. No. **22**, p. 19; English translation in *Dynamics of Star Clusters* (Eds. Goodman, J. and Hut, P.), D. Reidel Publ. Co., Holland, IAU Symposium No. 113, (1985), p. 521.
96. S. Chandrasekhar (1943), *Stochastic Problems in Physics and Astronomy*, Rev. Mod. Phys., **21**, 383.
97. S. Chandrasekhar (1943), *New methods in stellar dynamics* Ann. N. Y. Acad. Sci., **45**, p. 131.
98. S. Chandrasekhar and Von Neumann (1942), *The Statistics of the Gravitational Field Arising from a Random Distribution of Stars*, J. Astrophys., **95**, 489.
99. L. Landau (1936), *Physik. Zeits. Sowjetunion*, **10**, 154.
100. M. N. Rosenbluth, W.M. MacDonald and D.L. Judd (1957), *Fokker-Planck Equation for an Inverse-Square Force*, Phys. Rev., **107**, 1.
101. R. S. Cohen, L. Spitzer, Jr., and Paul McR. Routly (1950), *The Electrical Conductivity of an Ionized Gas*, Phys. Rev., **80**, 230.
102. E. M. Lifshitz and L. P. Pitaevskii (1981), *Physical Kinetics*, (Pergamon, London).
103. T. Padmanabhan (2000), *Theoretical Astrophysics: Volume 1, Astrophysical Processes*, Chapter 10, (Cambridge University Press, UK).
104. T. Padmanabhan (2009), *Random Walk Through Random Walks - II*, Resonance, **14**, 799.
105. G. N. Watson (1939) , *Three triple integrals*, Quarterly J. Math, **10**, 266.
106. M.L. Glasser and I.J. Zucker (1977), *Extended Watson integrals for the cubic lattice*, Proc. Natl. Acad. Sci.. USA, **74**, 1800.
107. P. G. Doyle and J. Laurie Snell (1984), *Random Walks and Electric Networks*, (Mathematical Association of America, Oberlin, OH).
108. P.L. Krapivsky and S. Redner (2004), *Random walk with shrinking steps*, Am. J. Phys., **72**, p. 591.

109. K. E. Morrison, *Random Walks with Decreasing Steps*, available at: <http://www.calpoly.edu/~kmorriso/Research/RandomWalks.pdf>
110. T. Padmanabhan (2008), *Thermodynamics of Self-Gravitating Particles*, Resonance, **13**, 941.
111. T. Padmanabhan (1990), *Statistical mechanics of gravitating systems*, Physics Reports **188**, 285.
112. S. Chandrasekhar (1939), *An Introduction to the Study of Stellar Structure*, (Dover).
113. T. Padmanabhan (1989), *Antonov instability and gravo-thermal catastrophe-revisited*, Astrophys. Jour. Supp. , **71**, 651.
114. V.A. Antonov (1962), Vest. Leningrad Univ. **7**, 135; English translation in *Dynamics of Star Clusters* (Eds. Goodman, J. and Hut, P.), D. Reidel Publ. Co., Holland, IAU Symposium No. 113, (1985), 525.
115. E.T. Whittaker (1927), *Electric Phenomena in a Gravitational Field*, Proc. Roy. Soc. Lond. A **116**, 720.
116. F. Rohrlich (1961), *The equations of motion of classical charges*, Ann. Phys., **13**, 93.
117. G. N. Plass (1961), *Classical electrodynamic equations of motion with radiative reactions*, Revs. Mod. Phys., **33**, 37.
118. H. Bondi and T. Gold (1955) , *The field of a uniformly accelerated charge, with special reference to the problem of gravitational acceleration* Proc. Roy. Soc. **A229**, 416.
119. M Born (1909), Ann. Physik., **30**, 1.
120. M. H. L. Pryce (1938), Proc. Roy. Soc., **A168**, 389
121. F. Rohrlich (1961), *The definition of electromagnetic radiation*, Nuovo Cimento, **21**, 811.
122. H. Padmanabhan and T. Padmanabhan (2010), *Aspects of electrostatics in a weak gravitational field*, Gen.Rel.Grav., **42**, 1153.
123. E. Eriksen and O. Gron (2004), *Electrodynamics of hyperbolically accelerated charges V. The field of a charge in the Rindler space and the Milne space*, Ann. Phys., **313**, 147.
124. E. Fermi (1921), Nuovo Cimento **22**, 176; reprinted in *Enrico Fermi, Collected papers (Note e memorie)* (Chicago University Press, Chicago, 1962); English translation in *Fermi and Astrophysics*, edited by V.G. Gurzadyan and R. Ruffini (World Scientific, Singapore, (2007)).
125. T. H. Boyer (1978), *Electrostatic potential energy leading to an inertial mass change for a system of two point charges*, Am. J. Phys., **46**, 383;
126. D. J. Griffiths and R. E. Owen (1983), *Mass renormalization in classical electrodynamics*, Am. J. Phys., **51**, 1120.
127. A. Gupta and T. Padmanabhan (1998), *Radiation from a charged particle and radiation reaction- revisited* Phys. Rev. D, **57**, 7241, [arXiv:physics/9710036],
128. P.A.M. Dirac (1938), *A New Basis for Cosmology*, Proc. Roy. Soc. **A165**, 199.
129. F. Rohrlich (1965), *Classical charged particles* (Addison - Wesley, Reading, MA).

Index

A

Abel's integral equation 103
accelerated charge
 field of 220
angular momentum 25, 43
 of electromagnetic field 241
Antonov instability 277

B

Bernoulli's equation 94, 96
black hole 3, 117, 127
 entropy 127, 133
 periodicity in imaginary time 195
 Schwarzschild metric 122
 surface gravity 132
 temperature 127, 194
 thermodynamics 128
blackbody radiation
 from black holes 127
 particles 232
 waves 232
Boltzmann equation 272
Born approximation 55
brachistochrone 73
 constant gravitational field 74
 history 77
 inverse square force 79

C

canonical ensemble 270
 phase transition 272
 specific heat 272
Cantor set 265
Casimir effect 209

 history 215
central force problem 28
central limit theorem 255
classical limit
 ray optics 10
 wave optics 10
Compton scattering 233
conic section 26, 39
constant gravitational field
 hodograph 76
 source for 58
Coriolis force 66, 143
 stability of motion 70
Coulomb field 28, 44, 49, 219
 mapping to oscillator 47
 accidental degeneracy 44
 energy levels 104
 motion in 39
 scattering cross section 54
Coulomb scattering
 curious features 52
cycloid 73

D

diffraction
 intuitive explanation 173
diffusion
 continuum random walk 249
 in velocity space 250
dimensional analysis 221
 electrostatics 110
 radiation field 221
dynamical friction 251
 Landau's derivation 253

E

- eccentricity 26
- electric field lines 279
 - in weak gravity 287
 - in gravitational field 279
 - radiation field 224
- electromagnetic field 231
 - angular momentum 241
 - blackbody cavity 231
 - collection of oscillators 206
 - energy-momentum tensor 235, 244
 - pair creation 199
 - particle motion 39
 - photons 231
 - relativistic definition 227
 - vector potential 14
- electrostatic potential
 - in weak gravity 284
 - weight of 290
- electrostatics 60
 - field of a line charge 109
 - fluid mechanics 90
 - Gauss law 90
 - infinities 110
- emergent gravity paradigm 134
- energy levels
 - hydrogen atom 43
 - potential 104
- Euclidean action
 - non-perturbative 199
 - Schwinger effect 198
 - tunneling 196
- Euclidean propagator
 - ground state energy 192
 - ground state wavefunction 192
 - thermal average 193

F

- Faraday's law 174
- Fokker-Planck equation 251
- Foucault pendulum 135
 - geometrical insight 146
 - rotation of the Earth 143
 - Thomas precession 150

G

- Galilean transformation 153
 - fluid flow 95
 - Schrödinger equation 156
- Gauss law
 - radiation 229

- general relativity 2, 117, 147
- golden ratio 265
- Green function 53

H

- Halley 77, 87
- Hamilton-Jacobi equation 46, 282
 - Coulomb problem 38
 - dispersion relation 13
 - electromagnetic field 14
 - from quantum theory 13
 - gravitational field 14
 - relativistic Coulomb problem 39
 - Schwarzschild metric 123
- Hamiltonian 9, 20, 45, 157
- harmonic oscillator 206
 - 4-dimensional 46
 - coherent states 162
 - shearing of 103
- Herman Melville 76
- hodograph 25
 - cycloidal motion 75
 - Kepler problem 25
- hydrogen atom
 - energy levels 43

I

- inverse square law 25, 27, 28, 32
 - source for 57
- inversion 60
- isochronous potential
 - definition 102
 - equidistant energy levels 106
 - quantum aspects 104
 - simple example 101
 - surprises in 101
- isothermal sphere 273
 - Antonov instability 277
 - non-singular solution 275
 - singular solution 274

J

- J.J. Thomson 220
- Jacobi-Mapertuis action 21, 181
 - area swept 35
 - relativistic propagator 186
 - tunneling 196

K

- Kepler problem 76
 - second focus 34
 - simple solution 25

Klein-Gordon equation 13, 159
 Lorentz invariance 159
 non-relativistic limit 159

L

Lagrange 65
 Lagrangian
 in rotating frame 66
 particle in a gravitational field 120
 Laplace equation
 fluid flow 89
 Larmor formula 234
 lattice of resistors 263
 latus rectum 26
 lemniscate 79
 lens equation 172
 Lobachevsky space 149
 Lorentz transformation 135
 properties 136

M

Maxwellian distribution 269
 Michael Grossmann 88
 microcanonical ensemble 270

N

Newton 34, 62, 77, 87
 non-relativistic limit
 of special relativity 153

O

optical system 170
 over-the-barrier-reflection 199
 complex path 201

P

parallel transport 147
 paraxial optics 168
 path integral propagator
 analytic continuation 191
 free particle 178
 non-relativistic 176
 paraxial optics 183
 quadratic actions 179
 relativistic 184
 Schrödinger equation 181
 square-root action 181
 stationary states 191
 transitivity constraint 177

Pauli matrices 50, 138
 period of oscillation
 1-dimensional motion 99
 potential 102
 scaling law 100
 phase space 16, 29
 Poisson equation 57, 60
 Green function 53
 inversion 60
 precession 37, 124
 Coulomb field 40
 elliptical orbits 37
 from Runge-Lenz vector 37
 general relativistic 125
 principle of equivalence 118

Q

quantum field theory 3
 the need for 189
 quantum gravity 3, 133
 quantum mechanics 2, 43
 wave-particle duality 2

R

radiation reaction 290
 rainbow 84
 secondary rainbow 86
 tertiary rainbow 87
 random walk 247
 continuum limit 248
 dimension dependence 261
 recurrent 259
 resistor network 264
 Rebecca Goldstein 87
 renormalization group 111
 running coupling constant 113
 Rindler frame 160, 280
 Runge-Lenz vector 45
 definition 31
 eigenvalues of 52
 quantum mechanics 44
 symmetry leading to conservation 36
 Rutherford scattering 27, 54

S

Schrödinger equation 7, 52, 177
 Coulomb problem 48
 delta function potential 111
 from Klein-Gordon equation 160
 in a constant field 158

- in a non-inertial frame 156
- Langer trick 105
- scattering 114
- Schwinger effect 196
- statistical mechanics
 - gravitating systems 269

T

- Thomas precession 135, 145
 - geometrical interpretation 135
 - intuitive interpretation 140
 - velocity space 148
- three-body problem
 - restricted 65
- Trojans 68

U

- uniformly accelerated observer 129

- temperature 131
- trajectory 129

V

- vacuum fluctuations 131, 132, 209
- vector potential 243
- Vena Contracta 97

W

- Watson integral 260
- wavefunction 7, 175
 - constructive interference 8
 - stationary phase 9
 - Wentzel-Kramers-Brillouin 15
- Wigner function 16