

Advanced Series in Mathematical Physics
Vol. 5

GEOMETRIC PHASES IN PHYSICS

Alfred Shapere

Frank Wilczek



World Scientific

Singapore • New Jersey • London • Hong Kong

Published by

World Scientific Publishing Co. Pte. Ltd.,
P O Box 128, Farrer Road, Singapore 9128
USA office: 687 Hartwell Street, Teaneck, NJ 07666
UK office: 73 Lynton Mead, Totteridge, London N20 8DH

The editors and publisher are grateful to the authors and the following publishers for their assistance and permission to reproduce the reprinted papers found in these volumes:

American Institute of Physics (*J. Chem. Phys.*);
American Physical Society (*Phys. Rev.*, *Phys. Rev. Lett.*);
Institute of Physics (*J. Phys. A*);
Kyoto University (*Prog. Theor. Phys.*),
Macmillan Magazines Ltd (*Nature*);
North-Holland Physics (*Nucl. Phys.*, *Phys. Lett.*);
Raman Research Institute (*Collected Works of S. Pancharatnam*);
Redakcja Acta Physica Polonica (*Acta Phys. Polonica*);
Springer-Verlag (*Commun. Math. Phys.*);
Taylor & Francis Ltd (*Mol. Phys.*);
The Faraday Society (*Disc. Farad. Soc.*);
The Royal Society (*Proc. Roy. Soc. London*).

While every effort has been made to contact the publishers of reprinted papers prior to publication, we have not been successful in a few cases. Where we could not contact the publishers, we have acknowledged the source of the material. Proper credit will be given to these publishers in future editions of this work after permission is granted.

Library of Congress Cataloging-in-Publication Data

Geometric phases in physics / edited by F. Wilczek & A. Shapere
p. cm. -- (Advanced series in mathematical physics; vol. 5)
1. Geometry. 2. Holonomy groups. 3. Mathematical physics.
I. Wilczek, Frank. II. Shapere, A. III. Series.
QC20.7.G44G46 1989 89-14624
530.1'5 -- dc20 ISBN 9971-50-599-1
ISBN 9971-50-621-1 (pbk)

Preface

During the last few years, considerable interest has been focused on a complex of physical ideas that share a common mathematical theme, the concept of holonomy. The recent flurry of activity began in 1984 with a paper by Michael Berry. He showed that the adiabatic evolution of energy eigenfunctions, with respect to a time-dependent quantum Hamiltonian $H(t)$, contains a phase of deeply geometrical origin (now known as “Berry’s phase”) in addition to the familiar dynamical phase

$$\exp - \frac{i}{\hbar} \int E(t) dt .$$

The additional phase approaches a finite, non-zero limit as the the Hamiltonian is taken more and more slowly around a closed path in its parameter space. Berry’s observation, although basically elementary, seems to be quite profound. Multiplicative phases—or more generally group transformations—with similar mathematical origins have been identified and found to be important in a startling variety of physical contexts, ranging from nuclear magnetic resonance to low Reynolds number hydrodynamics to quantum field theory. It now seems clear that Berry captured a particularly fruitful concept, of wide applicability.

There are several reasons for the impact of Berry’s work. Of course, the inherent universality and beauty of geometric phases has played a role, but it is also worth mentioning some of the extrinsic factors which made it “the right concept at the right time.” One factor was undoubtedly surprise—the surprise of the physics community that such a simple and fundamental aspect of the adiabatic theorem had been overlooked for so many years. Another reason for its impact was the emergence of gauge theories of the interactions of elementary particles. Many gauge theoretic ideas appear in the study of geometric phases, unencumbered by the complexities usually associated with relativistic quantum field theory. Conversely, ideas associated with geometric phases clarify some subtle issues in quantum field theory—as we shall find in Chapters 5 and 7. In an era of increasing separation between everyday reality and the more theoretical branches of physics, it has been refreshing and comforting to come across a concept that both helps to explain

some of the more abstruse ideas of quantum field theory, and leads to effects that can be readily measured in a laboratory.

Any judgement as to the value of an essentially mathematical concept, proposed for use as a tool in physics, should be at least partly based on its usefulness in practice. It is all too easy to believe that merely by adopting a new language one begins to make novel observations, but we trust that a perusal of the contents of this book will suffice to show that many genuinely new insights have been gained. Although it is at present used on a relatively modest scale, we believe that the concept of a geometric phase, repeating the history of the group concept, will eventually find so many realizations and applications in physics, that it will repay study for its own sake, and become part of the *lingua franca*.

The immediate origin of this book was a workshop held on “Non-integrable Phases in Dynamical Systems” in Minneapolis on 1-3 October 1987, under the aegis of the new Theoretical Physics Institute. The enthusiastic response of the participants, and the variety and quality of work presented, led us to think that a book on the subject would make a useful addition to the physics literature. The present volume contains reprints of papers we think are particularly important or instructive, together with several contributions which have not appeared in print before. The articles are arranged by subject in nine chapters, each of which begins with an introduction where we attempt to weave the varied material into a reasonably coherent whole.

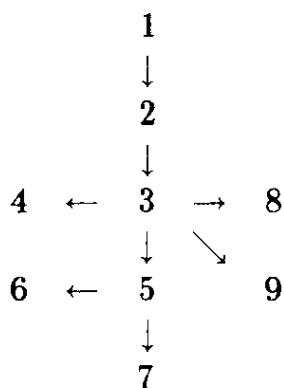
We do not purport to have made a comprehensive collection of relevant articles. Our choice of papers is more a function of our own particular interests and ignorance than anything else. Nevertheless, we hope that our book will serve as a useful introduction to an emerging field, and will stimulate its readers to seek out further material in their own areas of interest.

We are grateful to Michael Berry for his superb introductory survey and Daniel Arovas for his comprehensive article on fractional statistics and the fractional quantum Hall effect. We also wish to thank our editor, P.H. Tham, without whose hard work and assistance this book would not have been completed. The cover photo, “Calcutta Staircase 1988,” was taken by Catherine Shapere.

A Reader's Guide

Probably only the most adventurous readers will be motivated to read every chapter of this book, but many readers may be interested in browsing through unfamiliar territory. As an aid to field theorists who want to understand how non-Abelian gauge potentials apply to molecular systems, and to assist physical chemists in appreciating the connection between the molecular Aharonov-Bohm effect and gauge anomalies, we have prepared introductions to each of the chapters. The introductions serve several purposes. They provide some elementary background to topics covered in the chapter, touching at least briefly on each included article. They also try to put each chapter in perspective relative to the rest of the book, and to suggest further directions for research, when possible. It is our hope that by bringing together applications of geometric phases from a variety of fields, this book will inspire continued cross-fertilization between widely separated areas of physics.

Very roughly, the chapter dependence is as follows:



The first four chapters are of general interest. The first chapter includes two survey articles, by Berry [1.1]* and Jackiw [1.2], which we recommend to all readers. Both articles are relatively non-technical, and should help orient readers to the following three chapters. In addition, Berry's article includes original material on the natural metric on the projective Hilbert bundle and a detailed study of non-adiabatic corrections to the phase precession of the classical pendulum.

Chapter 2 contains some of the pre-1984 material which anticipated Berry's work on the quantal adiabatic phase, drawn from molecular physics and optics. This historical material is important and interesting in its own right—indeed, Pancharatnam's thirty-year-old paper on phase shifts of polarized light [2.1] has laid the groundwork for some modern optics experiments to measure geometric phases [4.2], and provides the basis for a recent extension of Berry's phase to non-closed paths [3.5]. The development over twenty years of phase concepts in molecular physics led to the use of gauge potentials in Born-Oppenheimer Hamiltonians several years before Berry's paper, and has borne a rich field of continuing activity.

* The notation $[M.N]$ refers to the N th article in chapter M .

The general foundations of our subject are laid out in Chapter 3. It contains Berry's original paper [3.1] and covers many subsequent extensions and elaborations of Berry's phase, such as Wilczek and Zee's non-abelian phase for degenerate Hamiltonians [3.3] and Aharonov and Anandan's phase for general cyclic evolution [3.4]. Its introduction includes, in an appendix, a pedagogical discussion of the mathematical context of Berry's phase, which may provide useful background for some of the more mathematical articles. The final article [3.7], which has not appeared previously, is a general discussion of the Born–Oppenheimer approximation and its field theory analogues, with phase effects and non-adiabatic corrections taken into account.

Basic applications of Berry's phase are treated in the following chapter, with articles drawn from optics, magnetic resonance, and molecular and atomic physics, from both the experimental and the theoretical literature. NMR and optics have provided some of the most successful tests of Berry's phase in macroscopic systems to date. In the fully quantum mechanical context of molecular physics, phase effects can lead to energy splittings and can shift quantum numbers, that have been observed experimentally.

The remaining five chapters are concerned with more specialized applications. Chapters 5, 6, and 7, respectively on fractional statistics, the quantized Hall effect, and anomalies and Wess–Zumino terms, are about geometric phases in many-body systems and quantum field theories. All three contain extensive introductions to aid the uninitiated. We would recommend reading Chapter 5 first, since the concept of fractional statistics plays a fundamental role in the theory of the fractional Hall effect, and is closely associated with Wess–Zumino terms. In fact, the boundaries between these three chapters are not too sharply defined—for instance, the review article by Arovas contains much general material on fractional statistics, although its main focus is on the fractional Hall effect. Also, these chapters touch on several topics which are not evident from their titles. The chapter on fractional statistics includes a paper by Laughlin on high-temperature superconductivity, and Chapter 6, on the quantized Hall effect, contains articles on two other two-dimensional systems—a network of current loops enclosing magnetic flux, and a Bloch electron in a transverse magnetic field.

Geometric phases also appear in classical systems. Hannay's angles are classical correspondants of Berry's phase, and can appear in any classical system described by action–angle variables, in response to adiabatic variation of the Hamiltonian. Another type of classical phase occurs in describing the motion of deformable bodies, which is especially useful in studying systems that are invariant under reparameterizations of time. In particular, the motion of a self-propelled body at low Reynolds number and the rotation of a self-deforming body in space can be described in terms of a gauge field over the space of shapes. These examples are all discussed in Chapter 8.

Finally, the last chapter contains Berry's elegant paper on higher-order corrections to the adiabatic approximation. It is our belief that there is

much room for further research in this area. Indeed, the subject of geometric phases in physics is far from closed, so perhaps it is fitting that we end on an unresolved note.

CONTENTS

Preface	v
A Reader's Guide	vii
Chapter 1 INTRODUCTION AND OVERVIEW	3
[1.1] M. V. Berry, "The Quantum Phase, Five Years After"*	7
[1.2] R. Jackiw, "Three Elaborations on Berry's Connection, Curvature and Phase," <i>Int. J. Mod. Phys. A3</i> (1988) 285–297	29
Chapter 2 ANTICIPATIONS	45
[2.1] S. Pancharatnam, "Generalized Theory of Interference, and its Applications," from Collected Works of S. Pancharatnam (Oxford University Press, UK, 1975)	51
[2.2] M. V. Berry, "The Adiabatic Phase and Pancharatnam's Phase for Polarized Light," <i>J. Mod. Optics</i> 34 (1987) 1401–1407	67
[2.3] G. Herzberg and H. C. Longuet-Higgins, "Intersection of Potential Energy Surfaces in Polyatomic Molecules," <i>Disc. Farad. Soc.</i> 35 (1963) 77–82	74
[2.4] A. J. Stone, "Spin-Orbit Coupling and the Intersection of Potential Energy Surfaces in Polyatomic Molecules," <i>Proc. R. Soc. Lond. A351</i> (1976) 141–150	80
[2.5] C. A. Mead and D. G. Truhlar, "On the Determination of Born-Oppenheimer Nuclear Motion Wave Functions Including Complications due to Conical Intersections and Identical Nuclei," <i>J. Chem. Phys.</i> 70 (05) (1979) 2284–2296	90
[2.6] Y. Aharonov and D. Bohm, "Significance of Electromagnetic Potentials in the Quantum Theory," <i>Phys. Rev.</i> 115 (1959) 485	104
Chapter 3 FOUNDATIONS	113
[3.1] M. V. Berry, "Quantal Phase Factors Accompanying Adiabatic Changes," <i>Proc. R. Lond. A392</i> (1984) 45–57	124

* Original Contribution.

[3.2]	B. Simon, "Holonomy, the Quantum Adiabatic Theorem, and Berry's Phase," <i>Phys. Rev. Lett.</i> 51 (1983) 2167–2170	137
[3.3]	F. Wilczek and A. Zee, "Appearance of Gauge Structure in Simple Dynamical Systems," <i>Phys. Rev. Lett.</i> 52 (1984) 2111–2114	141
[3.4]	Y. Aharonov and J. Anandan, "Phase Change during a Cyclic Quantum Evolution," <i>Phys. Rev. Lett.</i> 58 (1987) 1593–1596	145
[3.5]	J. Samuel and R. Bhandari, "General Setting for Berry's Phase," <i>Phys. Rev. Lett.</i> 60 (1988) 2339–2342	149
[3.6]	H. Kuratsuji and S. Iida, "Effective Action for Adiabatic Process," <i>Prog. Theo. Phys.</i> 74 (1985) 439–445	153
[3.7]	J. Moody, A. Shapere and F. Wilczek, "Adiabatic Effective Lagrangians"*	160

Chapter 4 SOME APPLICATIONS AND TESTS 187

[4.1]	A. Tomita and R. Chiao, "Observation of Berry's Topological Phase by Use of an Optical Fiber," <i>Phys. Rev. Lett.</i> 57 (1986) 937–940	193
[4.2]	M. V. Berry, "Interpreting the Anholonomy of Coiled Light," <i>Nature</i> 326 (1987) 277–278	197
[4.3]	J. Moody, A. Shapere and F. Wilczek, "Realizations of Magnetic-Monopole Gauge Fields: Diatoms and Spin Precession," <i>Phys. Rev. Lett.</i> 56 (1986) 893–896	199
[4.4]	D. Suter, G. C. Chingas, R. A. Harris and A. Pines, "Berry's Phase in Magnetic Resonance," <i>Mol. Phys.</i> 61 (1987) 1327–1340	203
[4.5]	R. Tycko, "Adiabatic Rotational Splittings and Berry's Phase in Nuclear Quadrupole Resonance," <i>Phys. Rev. Lett.</i> 58 (1987) 2281–2284	217
[4.6]	D. Suter, K. T. Mueller and A. Pines, "Study of the Aharonov-Anandan Quantum Phase by NMR Interferometry," <i>Phys. Rev. Lett.</i> 60 (1988) 1218–1220	221
[4.7]	A. Zee, "Non-Abelian Gauge Structure in Nuclear Quadrupole Resonance," <i>Phys. Rev.</i> A38 (1988) 1–6	224
[4.8]	C. A. Mead, "Molecular Kramers Degeneracy and Non-Abelian Adiabatic Phase Factors," <i>Phys. Rev. Lett.</i> 59 (1987) 161–164	230
[4.9]	B. Zygelman, "Appearance of Gauge Potentials in Atomic Collision Physics," <i>Phys. Lett.</i> A125 (1987) 476–481	234

* Original Contribution.

[4.10]	G. Delacrétaz, E. R. Grant, R. L. Whetten, L. Wöste and J. W. Zwanziger, "Fractional Quantization of Molecular Pseudorotation in Na_3 ," <i>Phys. Rev. Lett.</i> 56 (1986) 2598–2601	240
Chapter 5 FRACTIONAL STATISTICS		247
[5.1]	F. Wilczek and A. Zee, "Linking Numbers, Spin, and Statistics of Solitons," <i>Phys. Rev. Lett.</i> 51 (1983) 2250–2252	254
[5.2]	Y. S. Wu, "Multiparticle Quantum Mechanics Obeying Fractional Statistics," <i>Phys. Rev. Lett.</i> 53 (1984) 111–114	257
[5.3]	D. P. Arovas, R. Schrieffer, F. Wilczek and A. Zee, "Statistical Mechanics of Anyons," <i>Nucl. Phys.</i> B251 (1985) 117–126	261
Chapter 6 THE QUANTIZED HALL EFFECT		273
[6.1]	D. Arovas, J. R. Schrieffer and F. Wilczek, "Fractional Statistics and the Quantum Hall Effect," <i>Phys. Rev. Lett.</i> 53 (1984) 722–723	282
[6.2]	D. P. Arovas, "Topics in Fractional Statistics"*,	284
[6.3]	S. M. Girvin and A. H. MacDonald, "Off-Diagonal Long-Range Order, Oblique Confinement, and the Fractional Quantum Hall Effect," <i>Phys. Rev. Lett.</i> 58 (1987) 1252–1255	323
[6.4]	J. E. Avron, A. Raveh and B. Zur, "Quantum Conductance in Networks," <i>Phys. Rev. Lett.</i> 58 (1987) 2110–2113	327
[6.5]	R. B. Laughlin, "Superconducting Ground State of Noninteracting Particles Obeying Fractional Statistics," <i>Phys. Rev. Lett.</i> 60 (1988) 2677–2680	331
[6.6]	M. Wilkinson, "An Example of Phase Holonomy in WKB Theory," <i>J. Phys.</i> A17 (1984) 3459–3476	335
Chapter 7 WESS-ZUMINO TERMS AND ANOMALIES		355
[7.1]	M. Stone, "Born-Oppenheimer Approximation and the Origin of Wess-Zumino Terms: Some Quantum-Mechanical Examples," <i>Phys. Rev.</i> D33 (1986) 1191	361
[7.2]	J. Goldstone and F. Wilczek, "Fractional Quantum Numbers on Solitons," <i>Phys. Rev. Lett.</i> 47 (1981) 986	365
[7.3]	E. Witten, "Global Aspects of Current Algebra," <i>Nucl. Phys.</i> B223 (1983) 422	369

* Original Contribution.

[7.4]	I. J. R. Aitchison, "Berry Phases, Magnetic Monopoles, and Wess-Zumino Terms or How the Skyrmi got its Spin," <i>Acta Phys. Polonica</i> B18 (1987) 207–235	380
[7.5]	P. Nelson and L. Alvarez-Gaumé, "Hamiltonian Interpretation of Anomalies," <i>Commun. Math. Phys.</i> 99 (1985) 103–114	409
Chapter 8 CLASSICAL SYSTEMS		423
[8.1]	J. H. Hannay, "Angle Variable Holonomy in Adiabatic Excursion of an Integrable Hamiltonian," <i>J. Phys.</i> A18 (1985) 221–230	426
[8.2]	M. V. Berry, "Classical Adiabatic Angles and Quantal Adiabatic Phase," <i>J. Phys.</i> A18 (1985) 15–27	436
[8.3]	A. Shapere and F. Wilczek, "Gauge Kinematics of Deformable Bodies," to appear in <i>A. J. Phys.</i>	449
[8.4]	A. Shapere and F. Wilczek, "Geometry of Self-Propulsion at Low Reynolds Number," <i>J. Fluid Mech.</i> 198 (1989) 557–585	461
Chapter 9 ASYMPTOTICS		493
[9.1]	M. V. Berry, "Quantum Phase Corrections From Adiabatic Iteration," <i>Proc. R. Soc. Lond.</i> A414 (1987) 31–46	494

See also:

M.V. Berry & J.M. Robbins, "Geometric magnetism",
Proc Roy. Soc. Lond. A 442, 659 (95); 436, 631 (92).

GEOMETRIC PHASES IN PHYSICS

Chapter 1

INTRODUCTION AND OVERVIEW

[1.1]	M. V. Berry, "The Quantum Phase, Five Years After" [*]	7
[1.2]	R. Jackiw, "Three Elaborations on Berry's Connection, Curvature and Phase," <i>Int. J. Mod. Phys. A3</i> (1988) 285–297	29

* Original Contribution.

1

Introduction and Overview

The two articles contained in this chapter form the proper introduction to the main contents of the book. Here, we briefly discuss the general notion of holonomy, and illustrate it with a few primeval examples.

A phase is, for our purposes, not a state of matter but a complex number of unit modulus, an element of the group $U(1)$. We shall use the term somewhat loosely to encompass elements of matrix groups as well, such as $U(N)$. The phases we shall be interested in are often associated with cyclic evolution of a physical system. More specifically, we shall find that the cyclic variation of external parameters often leads to a net evolution involving a phase depending only on the *geometry* of the path traversed in parameter space. In other words, this phase is independent of how fast the various parts of the path are traversed. For non-cyclic evolution, the extra phase will depend on the endpoints of the path. The phase is non-integrable; it can not be written as a function just of the endpoints, because it depends on the geometry of the path connecting them as well.

The natural mathematical context for geometric phases is the theory of $U(N)$ fiber bundles. There one defines a phase, known as a holonomy,* that depends on the geometry of a loop, and is independent of any coordinate choice. (For a brief introduction to the world of fiber bundles, connections, and holonomy, see the introduction to Chapter 3.)

Examples of geometric phases abound in many areas of physics. Many familiar problems that we do not ordinarily associate with geometric phases may be phrased in terms of them. Often, the result is a clearer understanding of the structure of the problem, and an elegant expression of its solution.

Consider, for example, the precession of a Foucault pendulum. Standard treatments¹ calculate the rate of precession of a pendulum in a frame rotating with the surface of the earth in terms of the Coriolis force, but a much simpler and more geometric explanation may be given as follows. Suppose

* This phase is also, and perhaps more properly, known as an *anholonomy*. We shall adhere to the established mathematical terminology, although this other usage is quite common in the physics literature.

a pendulum is transported along a closed loop C , in the gravitational field of a point mass, and that the period and amplitude of its swing are small compared to the typical time and distance scales of the transport motion. We may assume that the loop lies on the surface of a sphere concentric with the mass, although this assumption is not necessary. Now when the pendulum returns to its initial position, its invariant plane will have rotated by some angle. For example, for transport around the sphere at a constant latitude of θ (relative to the north pole), a straightforward calculation shows that the net rotation will be $2\pi \cos \theta$ radians. A remarkable feature of this result is that it is independent of the rate at which different parts the loop are traversed (provided that the traversal is slow). This is a consequence of the fact that the Coriolis force is proportional to the velocity of transport, so that its integrated effect is invariant under rescalings of time. (Velocity-dependent forces, like the magnetic force on a moving charge, tend to be associated with geometric phases.) How does the pendulum precess as it is taken around a general path C ? For transport along the equator, the pendulum will not precess. This may be seen from a symmetry argument. The rate of precession does not depend on the direction of the pendulum's swing, so we may assume that the invariant plane lies in a north-south direction; then any precession would break the reflection symmetry between the northern and southern hemispheres, so the pendulum must not precess at all. (Alternatively, the Coriolis force at the equator always points vertically, and cannot torque the pendulum's invariant plane.) Now if C is made up of geodesic segments, the precession will all come from the angles where the segments meet; the total precession is equal to the net deficit angle, which in turn equals the solid angle enclosed by C modulo 2π . Finally, we can approximate any loop by a sequence of geodesic segments, so the most general result (on or off the surface of the sphere) is that the net precession is equal to the enclosed solid angle. This result may seem rather esoteric, but its generality and geometric nature suggest its depth. In fact, the mathematics describing it is essentially identical to that describing the motion of a charged particle in the field of a magnetic monopole, as well as interesting molecular, NMR, and optical systems.

A fundamental example of holonomy, involving a non-abelian symmetry, lies at the heart of general relativity. If a reference frame is parallel-transported around a closed loop in spacetime, then it is well known that the initial and final frames will not coincide. For causality's sake, we should really compare the result of parallel transport along two timelike curves with common endpoints. The final frames will be related by a Lorentz transformation, that is to say, an element of $SO(3, 1)$. The failure of the frames to coincide is an $SO(3, 1)$ holonomy. It may be used to measure the local curvature of spacetime. Thus, if we want to know the component $R_{\beta\mu\nu}^\alpha$ of the Riemann tensor, we should take a loop Γ in the $\mu\nu$ -plane, enclosing an infinitesimal area $dx^\mu dx^\nu$. The resulting holonomy M_β^α (an element of the Lie algebra of $SO(3, 1)$) from parallel transport around Γ will be proportional to

the area enclosed, to lowest order, and the constant of proportionality will be just the curvature component that we want:

$$M_\beta^\alpha = R_{\beta\mu\nu}^\alpha \cdot dx^\mu dx^\nu.$$

One last example concerns the motion of charged particles in strong magnetic fields. As is well known, in a constant and uniform magnetic field, a charged particle will move in a circle, or more precisely (in three dimensions) on a circular helix whose axis is parallel to the magnetic field direction. The motion is called cyclotron motion and the orbits cyclotron orbits, in reference to the use of such motion to guide particles at high-energy accelerators. A fundamental topic in plasma physics, with many applications in astrophysics, is how this motion is perturbed by various other effects such as inhomogeneities or time dependence in the magnetic field, electric fields, or gravitational forces. Insofar as the magnetic field is strong and reasonably homogeneous (*i.e.*, if it does not vary significantly over the radius of a cyclotron orbit), to a first approximation the motion is still cyclotron motion about the field lines, of gyrofrequency $\Omega(x) = eB(x)/mc$. The angular position of a particle relative to its guiding center axis will to lowest order be equal to the time integral of $\Omega(x(t))$. However, if the field lines are curved on a large scale, there will be both corrections to the angular position and drift corrections to the location of the guiding center.

Recently, Littlejohn² has introduced some fresh ideas and techniques, involving geometric phases, into this old subject. This has enabled him to derive the higher-order drifts much more easily and systematically than was previously possible. In fact, even without entering into any details one can see how geometric phases enter naturally into corrections to purely cyclotronic motion. As particles rotate about their guiding centers, they are “acquiring phase”—that is, their circular motion can be parametrized by an angle that increases with time. To a first approximation, the rate of phase accumulation is the local cyclotron frequency, uniquely determined by the local value of the magnetic field. However, if our particle returns to its starting point after following a closed field line, or after drifting through some loop, the total angle through which it has turned is not just the time integral of the local frequency, but contains an additional piece depending only on the geometry of the path through the space of possible magnetic field vectors. This additional phase, which is the leading correction to the total angle, is yet another example of holonomy in a purely classical system.

We hope that the above examples, which are concrete and readily visualized, will help to put in perspective some of the more abstract applications that follow.

[1] Keith R. Symon, *Mechanics*, 2nd edition (Reading: Addison-Wesley, 1960).

- [2] Robert G. Littlejohn, “Phase anholonomy in the classical adiabatic motion of charged particles,” *Phys. Rev.* **A38** (1988) 6034; “Geometry and guiding center motion,” in *Contemporary Mathematics*, edited by J.E. Marsden (Providence: American Mathematical Society, 1984), Vol. 28, p.151.

The Quantum Phase, Five Years After

M. V. Berry

H.H.Wills Physics Laboratory
Tyndall Avenue, Bristol BS8 1TL, U.K.

(Received 28 April 1988.)

ABSTRACT

Classical parallel transport of vectors is described in a manner immediately generalizable to parallel transport of quantum states in parameter space. The associated anholonomy is the geometric phase. One realization of parallel transport is by adiabatic cycling of the parameters. The phase is the flux of a 2-form. The 2-form is equivalent to the antisymmetric part of a gauge-invariant quantum geometric tensor. The symmetric part of this tensor gives a natural metric on parameter space. If the parameters are themselves regarded as dynamical variables, their adiabatic dynamics are influenced by a gauge field depending on both parts of the tensor. Corrections to the geometric phase (of higher order in an adiabatic parameter) can be obtained by successive transformations to moving frames, thereby generating a renormalization map of circuits in the space of Hamiltonians; the iterates diverge in a universal way. This quantum renormalization is illustrated by classical Newtonian and Hamiltonian renormalizations for a pendulum with changing frequency. To conclude, there are some historical remarks about geometric phases.

1. Introduction

The kind invitation to write this survey article provides two welcome opportunities. First, to present the fundamentals of the subject in a new perspective, reflecting some of the many recent developments and including some new material; and second, to make some historical remarks, drawing attention to important early works and describing the genesis of my own ideas in this field.

Two concepts are crucial to the understanding of this dusty corner of quantum theory which the brooms of our understanding are beginning to disturb. They are *anholonomy* and *adiabaticity*.

Anholonomy is a geometrical phenomenon in which nonintegrability causes some variables to fail to return to their original values when others, which drive them, are altered round a cycle. The simplest anholonomy is in the parallel transport of vectors, two examples being the change in the direction of swing of a Foucault pendulum after one rotation of the earth, and the change in the direction of linear polarization of light along a twisting ray [1][2] or coiled optical fibre [3-6] whose direction is altered in a cycle. The anholonomy to be described here is quantum-mechanical, and concerns the phase of a state which is parallel-transported round a cycle [7]. Parallel transport of a quantum state will here be introduced as a simple generalization of parallel transport of a vector.

Adiabaticity is slow change and therefore denotes phenomena at the border between dynamics and statics. Adiabatic change provides the simplest (but not the only [8]) way to make quantum parallel transport happen. The variables which are cycled are parameters in the Hamiltonian of a system. If the cycling is slow, the adiabatic theorem [9] guarantees that the system returns to its original state. But it usually acquires a nontrivial phase, a manifestation of anholonomy, and this is the phenomenon of interest here.

2. Classical Parallel Transport

It is convenient to begin by obtaining the law for the ordinary parallel transport of a vector over the surface of a sphere, expressing it in a form enabling instantaneous generalization to quantum mechanics. Let the unit vector \mathbf{e} be transported by changing the unit radius vector \mathbf{r} (Fig.1) and making two demands: that $\mathbf{e} \cdot \mathbf{r}$ must remain zero and that the orthogonal triad (frame) containing \mathbf{e} and \mathbf{r} must not twist about \mathbf{r} , i.e., $\boldsymbol{\Omega} \cdot \mathbf{r} = 0$ where $\boldsymbol{\Omega}$ is the angular velocity of the triad. These conditions define parallel transport of \mathbf{e} and lead to the law

$$\dot{\mathbf{e}} = \boldsymbol{\Omega} \wedge \mathbf{e} \quad \text{where} \quad \boldsymbol{\Omega} = \mathbf{r} \wedge \dot{\mathbf{r}} \quad (1)$$

This law is nonintegrable; when \mathbf{r} returns to its original direction after a circuit C on the sphere, \mathbf{e} does not return (in spite of never having been

twisted) but has turned through an angle $\alpha(C)$ which is the anholonomy now to be determined. Define $\mathbf{e}' \equiv \mathbf{r} \wedge \mathbf{e}$ (so that \mathbf{r} , \mathbf{e} , \mathbf{e}' form an orthogonal triad) and the complex unit vector

$$\psi \equiv (\mathbf{e} + i\mathbf{e}')/\sqrt{2} \quad (2)$$

in the plane perpendicular to \mathbf{r} . In terms of ψ , the parallel transport law (1) (which holds for \mathbf{e}' as well as \mathbf{e}) takes the simple form

$$\text{Im } \psi^* \cdot \dot{\psi} = 0 \quad \text{i.e.,} \quad \text{Im } \psi^* d\psi = 0 \quad (3)$$

where $d\psi$ is the change in ψ resulting from a change dr .

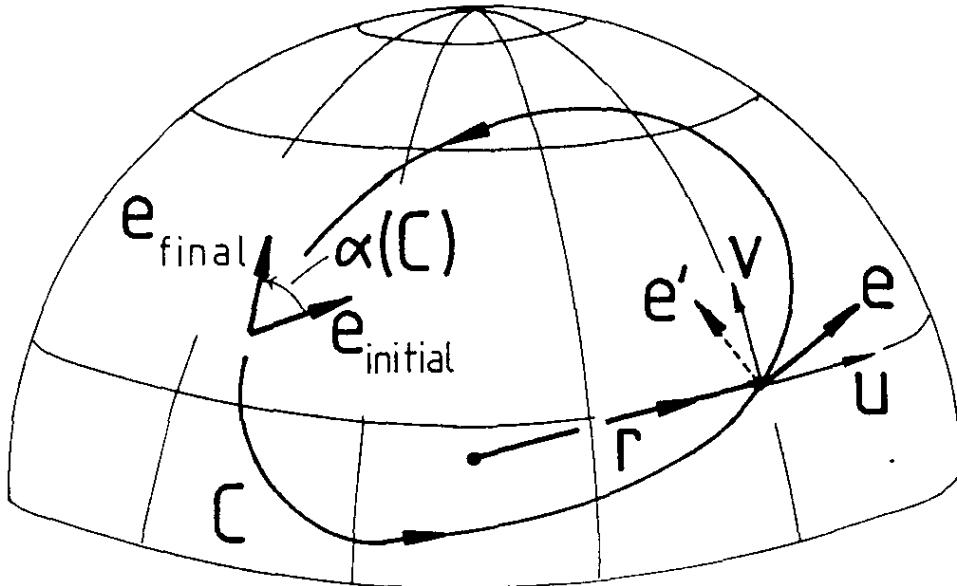


Figure 1. Rotation by $\alpha(C)$ after parallel transport of vector \mathbf{e} round circuit C on a sphere.

To find $\alpha(C)$ we chart the passage of \mathbf{e} and \mathbf{e}' relative to a local basis of unit vectors $\mathbf{u}(\mathbf{r}), \mathbf{v}(\mathbf{r})$ (Fig.1) defined at each point on the sphere. For example, we could choose \mathbf{u} and \mathbf{v} to lie along the parallel of latitude θ and meridian of longitude ϕ at $\mathbf{r} = (\sin\theta \cos\phi, \sin\theta \sin\phi, \cos\theta)$, i.e.,

$$\mathbf{u} = (-\sin\phi, \cos\phi, 0), \quad \mathbf{v} = (-\cos\theta \cos\phi, -\cos\theta \sin\phi, \sin\theta). \quad (4)$$

Specifying a local basis is equivalent to specifying the complex unit vector

$$\mathbf{n}(\mathbf{r}) \equiv (\mathbf{u}(\mathbf{r}) + i\mathbf{v}(\mathbf{r}))/\sqrt{2} \quad (5)$$

If the angle between the transported \mathbf{e} and the local \mathbf{u} is $\alpha(t)$, (2) and (5) give

$$\psi = \mathbf{n} \exp(-i\alpha) \quad (6)$$

whence (3) gives the anholonomy as

$$\begin{aligned} \alpha(C) &= \oint d\alpha = \text{Im} \oint \mathbf{n}^* \wedge \cdot d\mathbf{n} \\ &= \text{Im} \iint_{\partial S=C} d\mathbf{n}^* \cdot d\mathbf{n} \end{aligned} \quad (7)$$

where in the last equality Stokes' theorem has been used and the integral is over the area on the sphere bounded by C . It is an important result that the integrand in (7) is independent of the choice of local basis \mathbf{u}, \mathbf{v} : a change in this choice can be represented by a rotation $\mu(\mathbf{r})$ which induces the gauge transformation

$$\mathbf{n}(\mathbf{r}) \rightarrow \mathbf{n}'(\mathbf{r}) \exp\{i\mu(\mathbf{r})\} \quad (8)$$

under which $d\mathbf{n}^* \wedge \cdot d\mathbf{n}$ is invariant.

In terms of arbitrary parameters X_1, X_2 specifying \mathbf{r} (*i.e.*, position on the sphere), Eq. (7) can be written explicitly as

$$\alpha(C) = \text{Im} \iint_{\partial S=C} dX_1 dX_2 (\partial_1 \mathbf{n}^* \cdot \partial_2 \mathbf{n} - \partial_2 \mathbf{n}^* \cdot \partial_1 \mathbf{n}) \quad (9)$$

where ∂_j denotes $\partial/\partial X_j$. The choice $X_1 = \theta, X_2 = \phi$, together with (4), yields the integrand $d\theta d\phi \sin\theta$, which is simply the area element on the sphere, leading to the old result that the anholonomy $\alpha(C)$ is the *solid angle* subtended by C at the centre of the sphere.

3. Quantum Parallel Transport

To make the generalization to quantum mechanics, we replace the complex unit vector ψ by a normalized quantum state $|\psi\rangle$, *i.e.*, a unit vector in a Hilbert space, and position $\mathbf{r} = (X_1, X_2)$ on the sphere by position $X = (X_1, X_2, \dots)$ in a space of parameters governing the physical system represented by $|\psi\rangle$. At each X , $|\psi\rangle$ is defined up to a phase (just as \mathbf{e} was defined up to a rotation at each \mathbf{r}). Then a natural transport law [10] governing the phase of $|\psi\rangle$ as X varies is provided by reinterpreting (3) as the connection

$$\text{Im} \langle \psi | d\psi \rangle = 0. \quad (10)$$

Like (3), this law is nonintegrable: when X is taken round a circuit C , $|\psi\rangle$ returns with a changed phase. This change is the *quantum geometric phase* $\gamma(C)$; thus

$$\langle \psi_{\text{initial}} | \psi_{\text{final}} \rangle = \exp\{i\gamma(C)\}. \quad (11)$$

To find γ we again introduce a local basis by choosing at each X a definite (and so of course single-valued) state $|n(X)\rangle$, relative to which $|\psi\rangle$ is defined by

$$|\psi\rangle = |n(X)\rangle \exp(i\gamma) \quad (12)$$

Then (10) gives

$$\begin{aligned} \gamma(C) &= \oint d\gamma = -\text{Im} \oint \langle n|dn \rangle \\ &= -\text{Im} \iint_{\partial S=C} \langle dn| \wedge |dn \rangle \equiv - \iint_{\partial S=C} V(X). \end{aligned} \quad (13)$$

The integrand $V = \text{Im} \langle dn \wedge dn \rangle$ is the *phase 2-form*, whose flux through C gives the geometric phase. V is invariant under the gauge transformation

$$|n(X)\rangle \rightarrow |n'(X)\rangle \equiv |n(X)\rangle \exp\{i\mu(X)\} \quad (14)$$

For this mathematics to represent physics, it must be possible to implement the connection (10) by the Schrodinger equation

$$i\hbar|\dot{\Psi}\rangle = \hat{H}|\Psi\rangle \quad (15)$$

governing the evolution of any state $|\Psi\rangle$. A simple way [7] is to incorporate the parameters X into the Hamiltonian and change them slowly. Then the adiabatic theorem guarantees that in the absence of degeneracies (a restriction that can be removed [46]) $|\Psi\rangle$ will cling to one of the eigenstates of $\hat{H}(X(t))$, defined by

$$\hat{H}(X)|\psi\rangle = E_n(X)|\psi\rangle \quad (16)$$

The adiabatic ansatz

$$|\Psi\rangle \approx |\psi\rangle \exp \left\{ -\frac{i}{\hbar} \int_0^t dt' E_n(X(t')) \right\} \quad (17)$$

then gives the connection (10) immediately upon projecting (15) onto $|\psi\rangle$. The state $|n(X)\rangle$ in the 2-form (13) is any solution of (16) with a definite phase at each X .

Because of (17), the total phase change of $|\Psi\rangle$ includes a dynamical part as well as the $\gamma(C)$ being studied here. Thus

$$\langle \Psi_{\text{final}} | \Psi_{\text{initial}} \rangle = \exp\{i(\gamma_d + \gamma C)\} \quad (18)$$

where, for a circuit that takes a time T ,

$$\gamma_d = -\frac{1}{\hbar} \int_0^T dt E_n(X(t)) \quad (19)$$

One might say that γ_d and $\gamma(C)$ give the system's best answers to two questions about its adiabatic circuit. For γ_d the question is: how long did your journey take? For $\gamma(C)$ it is: where did you go?

Aharanov and Anandan [8] give a different interpretation of parallel transport. They regard the parameters X as labelling the state, rather than \hat{H} , so that X_1, X_2, \dots are coordinates in the *projective Hilbert space* that includes all quantum states, but where states differing only in phase (or normalization) are represented by the same point. Then a state $|\Psi\rangle$ evolving under Eq. (15) (not necessarily adiabatically) so as to return in T to the same X acquires a phase (18), with geometric part (13) (where the phase of $|n(X)\rangle$ is an arbitrary function of X) and dynamical part given by

$$\gamma_d = -\frac{1}{\hbar} \int_0^T dt \langle \Psi | \hat{H} | \Psi \rangle \quad (20)$$

instead of Eq. (19). The relation between the two approaches is that in the adiabatic case X parameterizes that part of the projective Hilbert space corresponding to the n th eigenstate of the chosen family of Hamiltonians $\hat{H}(X)$.

Several experiments have measured the geometric phase for particles, with spin 1/2 (neutrons [11]), spin 1 (photons [3]) and spin 3/2 (chlorine nuclei [12]). These depend on the result [7] that when \hat{H} is a rotationally symmetric function of the spin, *i.e.*,

$$\hat{H} = F(\boldsymbol{\sigma} \cdot \mathbf{X}) \quad (21)$$

where $\mathbf{X} = (X_1, X_2, X_3)$ and $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \sigma_3)$ is the vector spin operator, the geometric phase for the state with spin component n along \mathbf{X} is

$$\gamma_n(C) = -n \Omega(C) \quad (22)$$

where $\Omega(C)$ is the solid angle subtended by C at $\mathbf{X} = 0$.

These experiments all employ a superposition of eigenstates, rather than a single one, so that

$$\begin{aligned} |\Psi_{\text{initial}}\rangle &= \sum_n a_n |n\rangle \\ |\Psi_{\text{final}}\rangle &= \sum_n a_n |n\rangle \exp\{i(\gamma_{dn} + \gamma_n(C))\} \end{aligned} \quad (23)$$

At the end, that is after X has been cycled, an observable \hat{A} , which does not commute with the final \hat{H} , is measured (for example with a polarizer). Thus

$$\begin{aligned} \langle \hat{A} \rangle &= \sum_n |a_n|^2 \langle n | \hat{A} | n \rangle + 2 \operatorname{Re} \sum_{m \neq n} a_n^* a_m \langle m | \hat{A} | n \rangle \\ &\quad \times \cos \left\{ [\gamma_{dn} + \gamma_n(C)] - [\gamma_{dm} + \gamma_m(C)] \right\}. \end{aligned} \quad (24)$$

The oscillatory terms reveal $\gamma_n(C)$. This scheme has proved more convenient than the earlier suggestion [7] of splitting an ensemble of systems (*e.g.*, a beam of particles) into two subensembles, one being driven by an \hat{H} which is cycled and the other by an \hat{H} which is not, and then recombining the subensembles to detect $\gamma(C)$ by interference. (That is, instead of using one state and two Hamiltonians it is preferable to use two states — at least — and one Hamiltonian.)

Hannay [13] found an analogue of the geometric phase for *classical* systems. This was based on the simple observation that a quantum system in an eigenstate is an oscillator (because of the time factor $\exp(-iE_n t/\hbar)$), so that classical oscillators should exhibit similar anholonomy when parameters that govern them are cycled. The phase is now an angle, which may be an angle in space, like that of a wheel, or — more commonly — an abstract angle variable in phase space as with a harmonic oscillator. If the classical system is multiply periodic (integrable) for all X , with N freedoms (that is, coordinates $\mathbf{q} = (q_1 \cdots q_N)$ and momenta $\mathbf{p} = (p_1 \cdots p_N)$ and Hamiltonian $H(\mathbf{q}, \mathbf{p}; X)$, its orbit for fixed X winds round an N -torus [14] in phase space, with N angle variables $\boldsymbol{\theta} = (\theta_1 \cdots \theta_N)$ increasing uniformly. Conjugate to $\boldsymbol{\theta}$ are N adiabatically conserved actions $\mathbf{I} = (I_1 \cdots I_N)$ which label the torus. After a slow cycle of X the angles have acquired shifts which contain a geometric as well as a dynamical part. For a spinning particle [15-17] this classical anholonomy is the angle shift given by ordinary parallel transport of a vector.

Underlying Hannay's angles is a *classical 2-form*. This is the classical limit of the phase 2-form in Eq. (13), and semiclassical asymptotics [18] provides the expression

$$V(X) \xrightarrow[\hbar \rightarrow 0]{} -\langle d\mathbf{p} \wedge \cdot d\mathbf{q} \rangle / \hbar \quad (25)$$

whose symbols should be interpreted as follows. The wedge product \wedge links the d 's in parameter space. The scalar product \cdot links \mathbf{p} and \mathbf{q} . $\langle \cdot \rangle$ denotes an average over the angles on the torus labelled \mathbf{I} which at X corresponds [19] to the quantum state $|n\rangle$, *i.e.*, $\langle \cdot \rangle = \int_0^{2\pi} d\theta_1 \cdots d\theta_N / (2\pi)^N$. $d\mathbf{q}$ is the coordinate displacement linking corresponding points (labelled by the same $\boldsymbol{\theta}$) on the tori \mathbf{I} at X and $X + dX$, and similarly for $d\mathbf{p}$.

It is amusing to note that if the $2N$ variables \mathbf{q} and \mathbf{p} are replaced by the N complex variables

$$\mathbf{n} = (n_1 \cdots n_N) \equiv (\mathbf{q} + i\mathbf{p}) / \sqrt{2\hbar} \quad (26)$$

then (25) takes the form

$$V(X) \xrightarrow[\hbar \rightarrow 0]{} \text{Im} \langle d\mathbf{n}^* \wedge \cdot d\mathbf{n} \rangle \quad (27)$$

which bears a close formal resemblance both to the quantum expression (13) and the geometrical formula (7).

If the classical motion is not multiply periodic, that is if it is wholly or partly chaotic, the question of the classical limit of V is more delicate. It is tempting to claim that the limit is (25) for nonintegrable as well as integrable motion, but it is difficult to interpret the average $\langle \cdot \rangle$ and the displacements $d\mathbf{q}$ and $d\mathbf{p}$. In one of several interpretations, obtained by a semiclassical argument (not yet published) in collaboration with M. Wilkinson, $\langle \cdot \rangle$ denotes a time average over all points on an infinite orbit, and $d\mathbf{q}$ and $d\mathbf{p}$ link simultaneous points on the orbits for X and $X + dX$. For nonintegrable systems, however, it is not easy to express this result by replacing $\langle \cdot \rangle$ by a phase-space integral over the manifold explored by the orbit, because it is not clear what are then the ‘corresponding points’, linked by $d\mathbf{q}$ and $d\mathbf{p}$, on the manifolds for X and $X + dX$ (for an ergodic system these are the two constant-energy surfaces).

4. The Quantum Geometric Tensor

The central mathematical object underlying the quantum phase is the 2-form $V = \text{Im} \langle dn \wedge dn \rangle$. This is equivalent to an antisymmetric second-rank tensor field $V_{ij}(X)$ on the parameter space (or projective Hilbert space) with a quantum state $|n(X)\rangle$ defined at each point, namely

$$V_{ij}(X) = \text{Im}\{\langle \partial_i n | \partial_j n \rangle - \langle \partial_j n | \partial_i n \rangle\} \quad (28)$$

This tensor is invariant under the gauge transformation (14), but it is not the only such invariant tensor. More general is the *quantum geometric tensor*

$$T_{ij}(X) \equiv \langle \partial_i n | (1 - |n\rangle\langle n|) | \partial_j n \rangle \quad (29)$$

which is Hermitian, *i.e.*, $T_{ij} = T_{ji}^*$. The projector $|n\rangle\langle n|$ is essential to the gauge invariance. The imaginary part of T_{ij} is simply $V_{ij}/2$, so we can write

$$T_{ij} = g_{ij} + iV_{ij}/2 \quad (30)$$

where g_{ij} is the real symmetric tensor field $\text{Re } T_{ij}$.

We know the quantum meaning of V_{ij} : its flux gives the phase $\gamma(C)$. Therefore, it is natural to ask whether g_{ij} has significance. The answer is that g_{ij} provides a natural means of measuring distances along paths in parameter space; it is the *quantum metric tensor*. To understand why, observe that a natural measure of the squared distance between two nearby quantum states is the deviation from unity of their scalar product. If the states are $|1\rangle$ and $|2\rangle$ this gives, for the distance between the corresponding points X_1 and X_2 in parameter space,

$$\Delta s_{12}^2 = 1 - |\langle 1 | 2 \rangle|^2 \quad (31)$$

Taking the limit $1 \rightarrow 2$, and using the fact that all states are normalized, we obtain (using the summation convention for repeated indices i and j)

$$\begin{aligned} ds^2 &= \langle dn \cdot (1 - |n\rangle\langle n|) |dn \rangle = \langle \partial_i n | (1 - |n\rangle\langle n|) |\partial_j n \rangle dX_i dX_j \\ &= T_{ij} dX_i dX_j = g_{ij} dX_i dX_j, \end{aligned} \quad (32)$$

as claimed. The quantum tensor was introduced in an interesting paper by Provost and Vallee [50].

From its structure, g_{ij} can never give a negative ds^2 : in fact it is a positive semidefinite metric. Along a finite path (not necessarily closed) between $|1\rangle$ and $|2\rangle$, the quantum distance is

$$s_{12}(C) = \int_1^2 (g_{ij} dX_i dX_j)^{1/2}. \quad (33)$$

Page [33] and Bouchiat and Gibbons [41] give explicit forms for some metrics on the full Hilbert and projective Hilbert spaces.

The simplest example is a 2-state system, for which \hat{H} has the form (21), with $\hat{\sigma}$ the 3 Pauli matrices. If we take X as a unit vector, specified by parameters θ , ϕ (polar angles), the eigenstates are

$$|+\rangle = \begin{pmatrix} \cos(\theta/2) e^{i\phi/2} \\ \sin(\theta/2) e^{-i\phi/2} \end{pmatrix}, \quad |-\rangle = \begin{pmatrix} \sin(\theta/2) e^{i\phi/2} \\ -\cos(\theta/2) e^{-i\phi/2} \end{pmatrix} \quad (34)$$

For both of these, (32) gives $ds^2 = d\theta^2 + \sin\theta d\phi^2$, and this is the natural metric on the sphere of parameters (which in this case is also the projective Hilbert space).

Some interesting questions are suggested by this identification of g_{ij} as a metric on parameter space:

- (i) Do the geodesics, and in particular the shortest paths, connecting non-neighbouring states $|1\rangle$ and $|2\rangle$ have physical significance? One possibility, suggested by the work of Pancharatnam [20][21], is that the geodesics are the special paths along which the state preserves its phase in the sense that $\langle 1|2\rangle$ is real. This is true for the 2-state system just discussed, but seems to fail otherwise (probably for reasons of codimension). It is worth remarking that as $2 \rightarrow 1$ the overlap $\langle 1|2\rangle$ is real to second as well as first order in dX , for any path whatever.
- (ii) Can the geodesics be chaotic? This would require parameters X and states $|n(X)\rangle$ for which the Riemann curvature defined in terms of g_{ij} is negative (at least in some places) and the space is compact.
- (iii) Do *families* of geodesics (for example those issuing in different directions from the same point) exhibit the generic caustic singularities classified by catastrophe theory [22][23]? Do any such caustics have physical meaning? In 2-state systems the geodesics from X focus nongenerically at the

antipodal point on the sphere, where the state is orthogonal to $|n(X)\rangle$, but again this appears to be a special situation.

- (iv) Is there any meaning or interest in *quantizing* the geodesic motion in parameter space, for example by taking as Hamiltonian the Laplace-Beltrami operator $g^{-1/2}\partial_i g^{-1/2}g_{ij}\partial_j$ (where $g \equiv \det g_{ij}$)? Such quantizations are different from that described in the next section.

5. Dynamics of the Parameters

Until now we have regarded X as classical parameters which can be altered arbitrarily and which are unaffected by the quantum system they drive. But no physical action is unilateral and in reality X are themselves dynamical variables of a ‘heavy’ system coupled to the ‘light’ system (what we have so far called ‘the’ system) and therefore subject to reaction from it. Indeed the earliest application of the adiabatic theorem was the Born-Oppenheimer theory of molecules, in which X are coordinates describing the positions of the (heavy) nuclei and the light system is the electrons. Recently it has been pointed out [24–27] that in lowest order the reaction of the light system on the heavy dynamics is through a gauge field consisting of a vector potential whose curl is the phase 2-form V , and a scalar potential. Here I will show that what the gauge field really depends on is the quantum geometric tensor T_{ij} of section 3.

Let the heavy momenta, conjugate to X_i , be P_i . Then a fairly general nonrelativistic quantum Hamiltonian for the coupled system is

$$\hat{H}_{\text{tot}} = \frac{1}{2} \sum_{ij} Q_{ij} \hat{P}_i \hat{P}_j + H(\hat{\xi}; \hat{X}), \quad (35)$$

in which Q_{ij} is an inverse mass tensor, $\hat{\xi}$ are the dynamical variables of the light system (coordinates, momenta, spins, . . .) and H our previous Hamiltonian in which the X were regarded as parameters and which has eigenstates $|n(X)\rangle$ and energies $E_n(X)$. In the position representation for the heavy system, that is $\hat{P}_i = -i\hbar\partial_i$, the adiabatic ansatz is to write the full quantum state in the separated form

$$\langle X | \Psi \rangle \approx \Psi_{\text{heavy}}(X) | n(X) \rangle \quad (36)$$

and to consider the effective Hamiltonian governing Ψ_{heavy} to be

$$\hat{H}_{\text{eff}} = \langle n(X) | \hat{H}_{\text{tot}} | n(X) \rangle. \quad (37)$$

In \hat{H}_{eff} the reaction of the light on the heavy system comes from the action of the gradient operators \hat{P}_i on the X -dependence of $|n\rangle$. A straightforward calculation gives

$$\hat{H}_{\text{eff}} = \frac{1}{2} \sum_{ij} Q_{ij} \left\{ \hat{P}_i - A_i(\hat{X}) \right\} \left\{ \hat{P}_j - A_j(\hat{X}) \right\} + \Phi(\hat{X}) + E_n(\hat{X}) \quad (38)$$

where

$$A_i(X) = i\hbar \langle n | \partial_i n \rangle \quad (39)$$

and

$$\Phi(X) = \frac{\hbar^2}{2} \sum_{ij} Q_{ij} g_{ij}(X) \quad (40)$$

Here the emphasis is on the gauge potentials Φ and A_i — the scalar $E_n(X)$ is the ‘potential surface’ studied in conventional Born-Oppenheimer theory. Although (38) is a quantum Hamiltonian it can be used in suitable circumstances to calculate the *classical* motion of the heavy system, which will be affected by the fields A_i and Φ .

The physical effects of the vector potential A_i depend only on the ‘magnetic’ field

$$F_{ij} = \partial_i A_j - \partial_j A_i = -\hbar V_{ij} \quad (41)$$

(including its singularities and values in inaccessible regions — I am not denying the Aharonov-Bohm effect for heavy systems!). Thus the ‘magnetic’ field seen by the heavy system is the antisymmetric part of the quantum geometric tensor. The symmetric part of T_{ij} determines the ‘electric’ potential via Eq. (40). For an isotropic mass tensor, *i.e.* $Q_{ij} = \delta_{ij}/M$, Φ depends on $\text{Tr } g_{ij}$. It is a curious asymmetry that the ‘electric’ field depends on the gradients of g_{ij} , whereas the ‘magnetic’ field depends on V_{ij} itself.

The singularities of the gauge field are the *degeneracies* X^* of the spectrum, where $E_n(X^*) = E_{n\pm 1}(X^*)$. It is already known [7] that the ‘magnetic’ field V_{ij} (2-form) has monopole singularities. From the definition (29) of T_{ij} it is clear that g_{ij} has similar singularities, so that the ‘electric’ field near X^* is an inverse-cube force.

The situation near a degeneracy can be described by a special case of a simple model, which is of independent interest (and which has been studied from a different viewpoint by Anandan and Aharonov [28]), where the spin s of one (light) particle is coupled to the spatial coordinates of a second otherwise free (heavy) particle. Thus

$$\hat{H}_{\text{tot}} = \frac{1}{2M} \hat{P}^2 + F(\hat{\mathbf{X}} \cdot \hat{\boldsymbol{\sigma}}) \quad (42)$$

Near a degeneracy the appropriate model is a 2-state light system, so that we should take $s = \frac{1}{2}$, with linear coupling $F \propto \mathbf{X} \cdot \boldsymbol{\sigma}$.

The eigenvalues of $\mathbf{X} \cdot \hat{\boldsymbol{\sigma}}$ are nX , where $X \equiv |\mathbf{X}|$ and $-s \leq n \leq s$. The quantum tensor for the state $|n\rangle$ can be shown to be

$$T_{ij}^n(X) = \frac{1}{2X^2} \left\{ (s(s+1) - n^2) (\mathbf{e}_i \cdot \mathbf{e}_j - (\mathbf{e}_i \cdot \mathbf{x})(\mathbf{e}_j \cdot \mathbf{x})) \mp in(\mathbf{e}_i \wedge \mathbf{e}_j) \cdot \mathbf{x} \right\} \quad (43)$$

where $\mathbf{x} = \mathbf{X}/|\mathbf{X}|$ and \mathbf{e}_i is the unit vector along the i direction. The metric tensor g_{ij} has a zero eigenvalue, corresponding to radial parameter

displacements, which simply scale H leaving the states $|n\rangle$ unaffected: radial motions cover zero distance.

From Eqs. (38)–(40), the *classical* Newtonian equation for the heavy particle involves the Lorentz force from the magnetic monopole and the ‘electric’ force

$$-\nabla_{\mathbf{X}}\Phi(\mathbf{X}) = -\frac{\hbar^2}{2M}\nabla_{\mathbf{X}}\text{Tr}g_{ij} = \frac{\hbar^2(s(s+1)-n^2)}{MX^3}\mathbf{x} \quad (44)$$

This is of centrifugal type, and repels the parameters from a degeneracy (becoming significant at a distance of order $M^{-1/3}$), thereby tending to preserve the validity of the adiabatic approximation. We obtain, when the light particle is in the n th spin state,

$$M\ddot{\mathbf{X}} = \frac{S_z}{2X^3}\dot{\mathbf{X}}\wedge\mathbf{X} + \frac{(S^2-S_z^2)}{MX^4}\mathbf{X} - \frac{nF'(nX)}{X}\mathbf{X} \quad (45)$$

where $S_z \equiv n\hbar$ and $S^2 \equiv \hbar^2 s(s+1)$. This describes integrable motion, with conserved energy and modified angular momentum $M\mathbf{X}\wedge\dot{\mathbf{X}} - S_z\mathbf{X}/X$.

6. Adiabatic Renormalization

Now we return to the adiabatic scenario of section 3 and realize that γ_d and $\gamma(C)$ in Eq. (18) are but the first two terms in an infinite series involving powers of an adiabatic slowness parameter ϵ , influencing the dynamics through \hat{H} whose time-dependence enters in the combination ϵt . The dominant term is γ_d (Eq. 19) and is of order ϵ^{-1} . The next term is $\gamma(C)$, whose unique feature — and the reason for its being called geometric — is that it is independent of ϵ , and so depends only on the sequence of Hamiltonians along the circuit and not on its time history.

This uniqueness is not threatened by the observation that transformation to a moving frame (a common practice in problems involving spin [11]) can make $\gamma(C)$ appear ‘dynamical’ by making it emerge from a correction to the energy rather than as anholonomy: the geometric structure of $\gamma(C)$ is independent of how it is derived.

Transformations to moving frames have however another interest, in that they form the basis of a renormalization (iteration) technique for generating higher-order corrections to the phase. Details of the technique have been published elsewhere [29]; here I will outline the central idea, and give an example.

Let the Hamiltonian $\hat{H}_0(t)$ generating the quantum motion be cyclic, in the sense that $\hat{H}_0(+\infty) = \hat{H}_0(-\infty)$, and let it have instantaneous eigenstates $|n_0(t)\rangle$ and energies $E_0(n,t)$. The evolving state $|\Psi_0(t)\rangle$ is determined by

$$i|\dot{\Psi}_0(t)\rangle = \hat{H}_0(t)|\Psi_0(t)\rangle \quad (46)$$

with the initial condition

$$|\Psi_0(-\infty)\rangle = |n_0(-\infty)\rangle \equiv |N\rangle \quad (47)$$

After the cycle, *i.e.*, at $t = +\infty$, $|\Psi_0\rangle$ will have returned only approximately to $|N\rangle$, so a phase can be defined precisely by

$$\gamma \equiv \text{Im} \log \langle N|\Psi_0(+\infty)\rangle - \gamma_d \quad (48)$$

The geometric phase $\gamma(C)$ (Eq. 13) is $\lim_{\epsilon \rightarrow 0} \gamma$. The aim is to obtain increasingly accurate approximations to $\gamma - \gamma(C)$. It is worth emphasizing that the non-aim is the determination of the nonadiabatic transition probability $1 - |\langle N|\Psi_0(+\infty)\rangle|^2$, because this is the usual objective of adiabatic theory, and that the non-method is perturbation theory, because this is the usual technique [30][49].

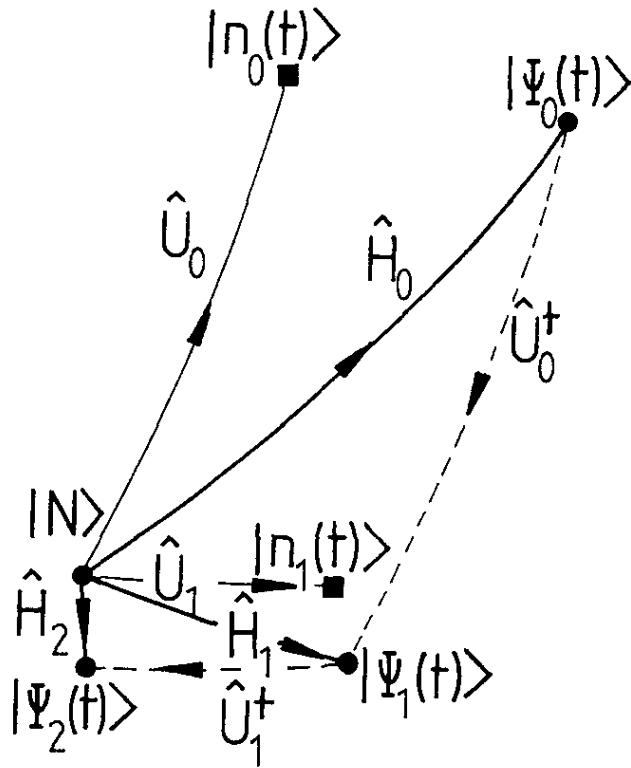


Figure 2. Renormalization in Hilbert space.

To explain the method used instead, we refer to Fig. 2. When ϵ is small we expect $|\Psi_0(t)\rangle$ to be close to $|n_0(t)\rangle$. This suggests that defining a unitary transformation $\hat{U}_0(t)$ by

$$|n_0(t)\rangle = \hat{U}_0(t)|N\rangle \quad (49)$$

will be useful. The inverse operator \hat{U}_0^\dagger sends $|n_0(t)\rangle$ back to $|N\rangle$, that is, it freezes the moving eigenstate. Therefore \hat{U}_0^\dagger should almost freeze the evolving state $|\Psi_0(t)\rangle$, and so we define

$$|\Psi_1(t)\rangle \equiv \hat{U}_0^\dagger|\Psi_0(t)\rangle. \quad (50)$$

We are attempting to follow $|\Psi_0(t)\rangle$ by transforming to a moving frame. The Hamiltonian governing $|\Psi_1\rangle$ is

$$\hat{H}_1 = \hat{U}_0^\dagger \hat{H}_0 \hat{U}_0 - i \hat{U}_0^\dagger \dot{\hat{U}}_0. \quad (51)$$

in which the second term is the quantum analogue of the inertial forces generated classically by transforming to a moving frame.

Now that the original problem has been reduced to one of the same form but involving $|\Psi_1\rangle$ and \hat{H}_1 instead of $|\Psi_0\rangle$ and \hat{H}_0 , it is natural to iterate the process by defining $|\Psi_2\rangle \equiv \hat{U}_1^\dagger |\Psi_1\rangle$, where \hat{U}_1^\dagger freezes the eigenstates $|n_1\rangle$ of \hat{H}_1 . This defines a *renormalization map* $\hat{H}_k \rightarrow \hat{H}_{k+1}$ in Hamiltonian space. The form of the map is simple when written in a basis of initial states (which are unaffected by renormalization) and with the phases of the eigenstates chosen so that they are parallel-transported, *i.e.*, $\langle n_k | \dot{n}_k \rangle = 0$:

$$\langle M | \hat{H}_{k+1} | N \rangle = E_k(n, t) \delta_{MN} - i \frac{\langle m_k(t) | \dot{H}_k(t) | n_k(t) \rangle}{E_k(m, t) - E_k(n, t)} (1 - \delta_{MN}) \quad (52)$$

The k th approximant $\gamma^{(k)}$ to the phase is obtained by neglecting the off-diagonal terms in \hat{H}_{k+1} . $\gamma^{(k)}$ is the sum of the phase anholonomies of the Hamiltonians $\hat{H}_0 \dots \hat{H}_k$ (arising from the continuation of $n_k(t)$) from $t = -\infty$ to $t = +\infty$ and reflected as phase factors $\langle N | U_k(+\infty) | N \rangle$, together with an additional term involving E_k [29]. (A contrary choice of phases, *i.e.*, $|n_k(+\infty)\rangle = |n_k(-\infty)\rangle$, gives $\langle N | U_k(+\infty) | N \rangle = 1$, but now the diagonal terms in Eq. (52) contain extra terms $-i\langle n_k | \dot{n}_k \rangle$ and all corrections — including $\gamma^{(0)} = \gamma(C)$ as mentioned previously — appear dynamical.)

Each renormalization produces a new Hamiltonian which over $-\infty < t < +\infty$ traverses a loop in Hamiltonian space. If the renormalizations converged, successive loops would get smaller (by a factor ϵ each time). But this does not, and indeed cannot, happen. If it did, $\langle \Psi(-\infty) | \Psi(+\infty) \rangle$ would have modulus unity, contradicting the existence of transitions to other states. The accumulation of inertial forces in successive renormalizations defeats our attempts to follow the motion, which slips out of control, causing the scheme to diverge.

Nevertheless, the corrections generated by renormalization do get smaller at first, and enable γ to be determined with an error of order $\exp(-1/\epsilon)$, which occurs after $k \sim 1/\epsilon$ renormalizations. A detailed exploration [29] of 2-state systems (the simplest nontrivial case, for which the geometry of the loop map can be made explicit) reveals that the Hamiltonian loops (which lie on a 2-sphere) get smaller and then larger in a universal way (that is, almost always independent of the form of the initial loop).

This procedure is typical of asymptotic procedures and occurs also in the more usual adiabatic perturbation theory. It prompts interesting questions. What is the dynamical significance of the moving frame that produces the best approximant to γ , generated by $\hat{U}_{k \sim 1/\epsilon} \hat{U}_{k-1} \cdots \hat{U}_0$? Can the exponential

residue $\gamma - \gamma^{(k)}$ be more closely approximated by generalizing the Borel (or some other) resummation method [47]?

It is instructive to illustrate adiabatic renormalization with the *classical* problem which gave birth to the entire subject, namely the Ehrenfest-Einstein pendulum [31] whose frequency is slowly changed. Newton's equation is

$$\partial_t^2 x(t) + \omega^2(t) x(t) = 0 \quad (53)$$

in which the frequency $\omega(t)$ is a smooth nonzero function with $\omega(+\infty) = \omega(-\infty) \equiv \omega_\infty$. The same equation describes the (time-independent) quantum mechanics of a beam of particles with energy E encountering a potential well or hill $V(x)$ such that $E > V(x)$ for all x .

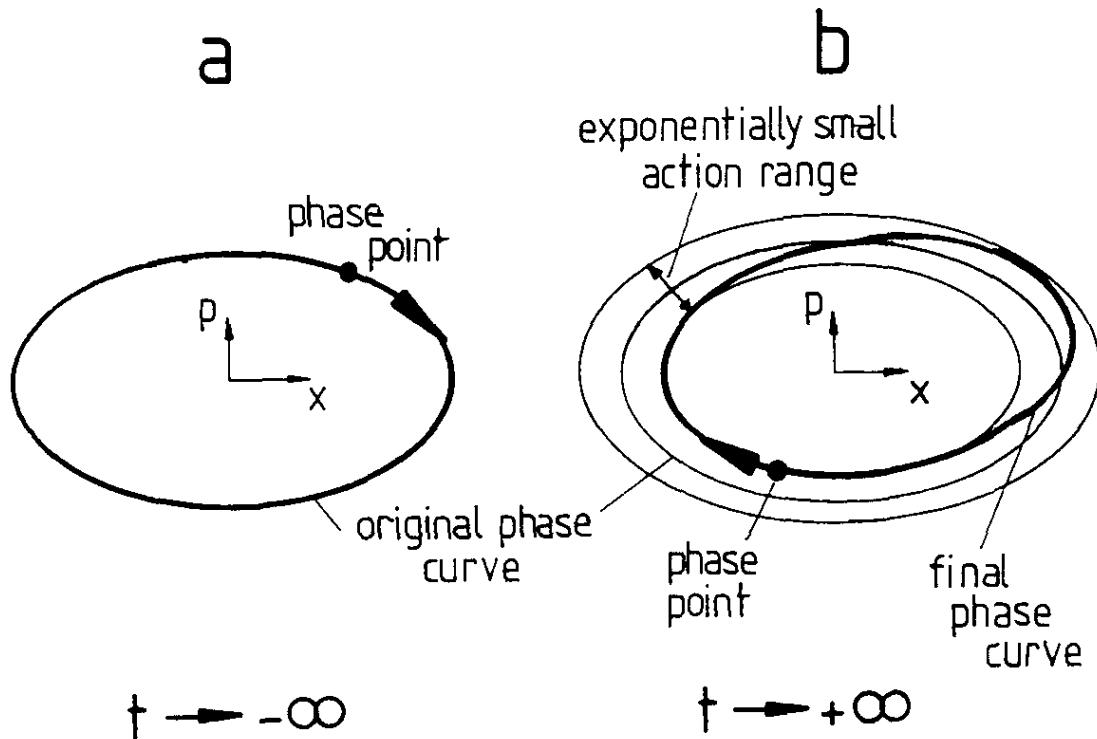


Figure 3. (a) Initial and (b) final phase portraits for slowly-altered pendulum.

Consider motion in the phase plane with variables x and $p = \dot{x}$. Initially, i.e., as $t \rightarrow -\infty$, each phase point moves round an ellipse with frequency ω_∞ (Fig. 3a). The subsequent motion lies on a curve that at each instant approximates one of the elliptical contours of the Hamiltonian

$$H(x, p, t) = \frac{1}{2} [p^2 + \omega^2(t)x^2] \quad (54)$$

at that time. These subsequent ellipses have approximately the same area as the original one, because the adiabatically-conserved action is $\text{area}/2\pi$.

As $t \rightarrow +\infty$, then, the phase point is close to its original ellipse and we can ask: where is it on the ellipse, i.e., what is its phase?

This would be a question about Hannay's angle were it not for the fact that this is a classical problem without anholonomy, so that the angle we seek consists entirely of nonadiabatic corrections. By identifying the solution of (53) with quantum transmitted and reflected waves, it can be shown that the oscillation which begins as

$$x + ip/\omega_\infty = A \exp \left\{ -i \left(\int_0^t dt' \omega(t') + \sigma \right) \right\} \quad (t \rightarrow -\infty) \quad (55)$$

ends as

$$x + ip/\omega_\infty = A \exp \left\{ -i \left(\int_0^t dt' \omega(t') + \sigma \right) \right\} [T^{-1} + R T^{-1} \exp(2i\sigma)] \quad (t \rightarrow +\infty) \quad (56)$$

where A is a real constant and R and T are the complex quantal reflection and transmission coefficients.

Therefore the phase shift depends on the initial phase σ , but this dependence is slight because R is exponentially small in the slowness parameter ϵ (if ω depends on ϵt). In any case, we can define a phase by averaging over σ , with the exact result

$$\gamma \equiv -\frac{1}{2\pi} \int_{-\pi}^{\pi} d\sigma \lim_{t \rightarrow \infty} \left[\text{Im} \log (x + ip/\omega_\infty) + \int_0^t dt' \omega(t') \right] = \text{Im} \log T. \quad (57)$$

Thus 'Hannay's angle' is here the phase of the transmission coefficient. The final action I also depends on σ , but the range is $(I_{\max} - I_{\min})/I_{\text{initial}} = 4|R|/|T|^2$ which again is of order $\exp(-1/\epsilon)$; the whole initial ellipse of phase points evolves ultimately into one exponentially close to it and deforming periodically with frequency ω_∞ (Fig. 3b). Newtonian renormalization of (53) is based on the transformation

$$x(t) \equiv \frac{x_1(t_1)}{\omega^{1/2}(t)}; \quad t_1 \equiv \int_0^t dt' \omega(t') \quad (58)$$

whose new coordinate satisfies

$$\partial_{t_1}^2 x_1(t_1) + \omega_1^2(t_1) x_1(t_1) = 0 \quad (59)$$

where

$$\omega_1^2(t_1(t)) = 1 + \omega^{-3/2} \partial_t^2 \omega^{-1/2}. \quad (60)$$

Clearly $\omega_1 \approx 1$ if ω varies slowly.

Renormalization consists of iterating this transformation, the aim being to freeze the frequency. The k th approximant for γ is obtained by approximating $\omega_{k+1} \approx 1$, so

$$\gamma^{(k)} = \int_{-\infty}^{\infty} (dt_{k+1} - \omega(t) dt) = \int_{-\infty}^{\infty} dt \omega \left(\prod_{j=1}^k \omega_j - 1 \right) \quad (61)$$

Thus

$$\begin{aligned} \gamma^{(0)} &= 0 && \text{(no anholonomy in this problem)} \\ \gamma^{(1)} &= \int_{-\infty}^{\infty} dt \omega(t) (\omega_1(t_1(t)) - 1) \end{aligned} \quad (62)$$

etc. $\gamma^{(1)}$ is of order ϵ .

An equivalent *Hamiltonian* renormalization is produced by iteration of the canonical transformation generated by

$$S(x, p_1, t) = x p_1 \omega^{1/2} - x^2 \partial_t \omega / 4\omega. \quad (63)$$

This gives

$$x_1 = x \omega^{1/2}; \quad p_1 = p \omega^{-1/2} + x \partial_t \omega / 2\omega^{3/2} \quad (64)$$

and hence the transformed Hamiltonian

$$\bar{H}(x_1, p_1, t) = H + \partial_t S = \frac{1}{2} \omega(t) (p_1^2 + \omega_1^2 x_1^2) \quad (65)$$

where ω_1 is given by (60). Rescaling time to t_1 as defined in (58) now gives

$$H_1(x_1, p_1, t_1) = \frac{1}{2} (p_1^2 + \omega_1^2(t_1) x_1^2(t_1)) \quad (66)$$

which is the first renormalization of the original Hamiltonian (54). The aim of subsequent renormalizations is to freeze the Hamiltonian into one whose contours are circles.

I have expressed these classical iteration schemes in terms of the renormalization of Newton's or Hamilton's equations in order to illustrate the idea behind the quantum renormalization described earlier. But they can be shown to be equivalent to the following fairly conventional WKB-like [32] procedure (to be contrasted with an unconventional WKB analysis by Wilkinson [45] which, unlike this one, does involve anholonomy). Write the exact solution of Eq. (53) as

$$x(t) = \Omega^{-1/2}(t) \cdot \exp \left\{ i \int_0^t dt' \Omega(t') \right\}. \quad (67)$$

Then the 'frequency' $\Omega(t)$ satisfies

$$\Omega^2(t) = \omega^2(t) + \Omega^{1/2}(t) \partial_t^2 \Omega^{-1/2}(t). \quad (68)$$

In terms of Ω , the phase shift is, exactly,

$$\gamma = \int_{-\infty}^{\infty} dt [\Omega(t) - \omega(t)]. \quad (69)$$

Successive approximants are obtained by the iteration

$$\Omega^{(0)} = \omega; \quad \Omega^{(k+1)} = \left[\omega^2 + (\Omega^{(k)})^{1/2} \partial_t^2 (\Omega^{(k)})^{-1/2} \right]^{1/2}. \quad (70)$$

The inevitability and universality of the divergence of these schemes can be demonstrated by considering high-order iterations of Eq. (60), for which

$$\omega_k(t) \equiv 1 - \delta_k(t) \quad (71)$$

and $\delta_k \ll 1$. Then $t_{k+1} \approx t_k$ (cf.Eq. (58)), and Eq. (60) can be written approximately as

$$\delta_{k+1}(t) \approx -\frac{1}{4} \partial_t^2 \delta_k(t). \quad (72)$$

The asymptotics of this recursion as $k \rightarrow \infty$ can be estimated by Fourier analysis, on the assumption that $\delta_0(t)$ is a real function of $\tau \equiv \epsilon t$, analytic in a strip about the real τ axis with its nearest singularities at $\tau_1 \pm i\tau_2$. Then with $\xi \equiv (\epsilon t - \tau_1)/\tau_2$ it is possible to show that

$$\delta_k(t) \xrightarrow[k \rightarrow \infty]{} \left[\frac{A \epsilon^{2k} (2k)!}{4^k \tau_2^{2k+1}} \right] \left[\frac{\cos \{(2k+1) \cos^{-1}(1+\xi^2)^{-1/2}\}}{(1+\xi^2)^{k+1/2}} \right] \quad (73)$$

where A is a constant.

The first factor in (73) shows the divergence: $\epsilon^{2k} (2k)!$ decreases until $k \sim \tau_2/\epsilon$, when $\delta_k \sim \exp(-2\tau_2/\epsilon)$, and then increases until $\delta_k \sim 1$, when the scheme breaks down. The second factor is the universal function describing the asymptotic ‘frequency.’

7. Historical Remarks

First I consider the important special case where the transported states $|\psi\rangle$ can be represented by wavefunctions that are *real*. Then the only possible phase factors associated with a circuit C are ± 1 . It follows [7] from the result (22) for spins that the factor is -1 when C encloses a degeneracy X^* of the spectrum to which $|\psi\rangle$ belongs; otherwise, it is +1. The peculiarity of this case is that parallel transport (10) is the only possible smooth continuation law, rather than a mathematically natural choice, concordant with quantum dynamics, from a infinity of possibilities.

Eigenfunctions can always be made real if their Hamiltonian matrix is real symmetric rather than complex Hermitian (this is the case when there is (bosonic) time-reversal symmetry [34]). Thus the phase law states

that an eigenfunction of a real symmetric matrix depending on parameters *changes sign* under smooth continuation round a degeneracy. This result is so simple - it holds even for 2×2 matrices - as to deserve mention in elementary expositions of matrix theory, but I have not found it in any such text. Arnold [14] is aware of the sign change, and attributes it to Uhlenbeck [35] in 1976. It was already known to theoretical chemists: Herzberg and Longuet-Higgins [36] gave an explicit statement in 1963. But the sign change (for 2×2 matrices) was implicit in work of Darboux [37] as long ago as 1896. This concerns the differential geometry of surfaces, and is worth describing.

Darboux considered a curved surface described locally by its deviation $z(X_1, X_2)$ from the plane $X = (X_1, X_2)$. Then the 2×2 real symmetric curvature matrix at X is

$$H_{ij}(X) = \partial_i \partial_j z(X). \quad (74)$$

The two eigenvalues are the principal curvatures at X , and the corresponding eigenvectors give the (orthogonal) directions of the lines of curvature at X . Degeneracies are *umbilic points*, where the surface is locally spherical (two curvatures equal). Umbilics are singularities of the net of curvature lines. The sign-change rule states that a line of curvature turns by π in a circuit of an umbilic: the Poincaré index of the tensor field (74) is $\pm \frac{1}{2}$. Fig. 4 shows how this happens for the three generic patterns [38][39] of curvature lines near an umbilic; the star has index $-\frac{1}{2}$, and the lemon and monstar have index $+\frac{1}{2}$. Star and lemon singularities occur as disclinations in liquid crystals [48].

The full phase — rather than the impoverished special case of the sign change for real matrices — was anticipated at least twice. First, in the mid-1950's, Pancharatnam [20][21][40] studied the 2-state Hermitian case in the context of the polarization states of light travelling in a fixed direction. The parameter space is the surface of the Poincaré sphere. Pancharatnam introduced the useful idea of defining two different states $|1\rangle$ and $|2\rangle$ as ‘in phase’ if the intensity of their superposition is maximal, a condition equivalent to their overlap $\langle 1|2\rangle$ being real and positive. This defines a connection between the corresponding parameters X_1 and X_2 as the state $|2\rangle$ obtained from $|1\rangle$ by phase-preserving transport along the shorter geodesic arc between X_1 and X_2 . He discovered that the connection is nontransitive: a circuit $X_1 X_2 X_3 X_1$ produces a state differing from $|1\rangle$ by precisely the same phase anholonomy [21] (minus half the solid angle of the circuit) as that given by parallel transport.

Second Mead [24] and Mead and Truhlar [42], studying adiabatic theory for molecules, made two important advances. They showed how the sign-change rule for degeneracies would induce modifications in the nuclear dynamics and hence change the vibration-rotation spectrum. And they realized that in the absence of time-reversal symmetry the nuclear dynamics would be influenced by the vector potential (39) and the corresponding ‘magnetic’ field (41), for which they gave a general formula.

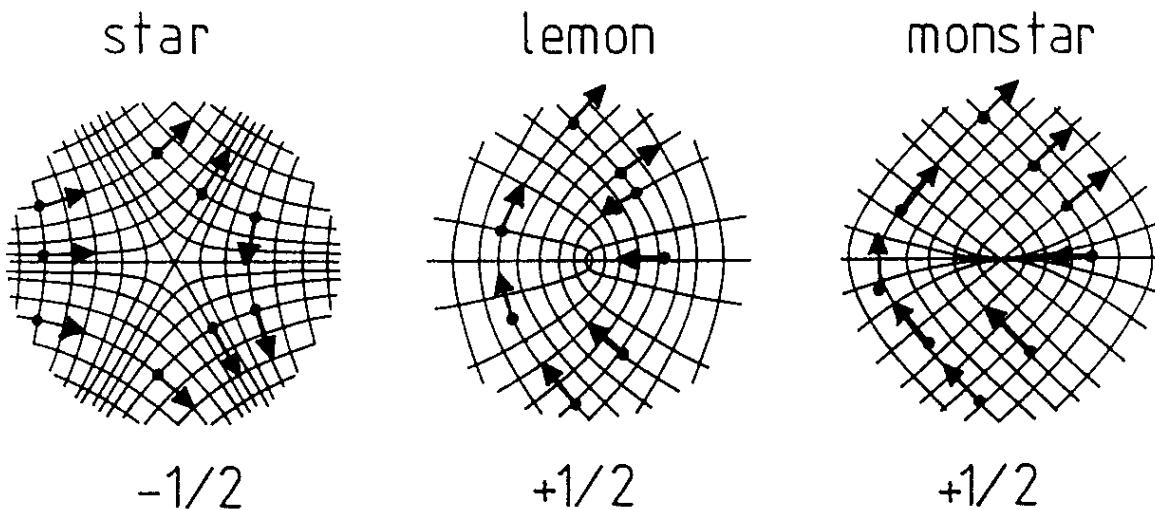


Figure 4. The generic patterns of lines of curvature near an umbilic point on a surface, illustrating the reversal ($\pm 1/2$ index) round the singularity.

My involvement with this subject began in 1979 with the appreciation [43] that degeneracies play a part in determining the fine-scale statistics of energy levels of quantum systems whose classical counterparts are nonintegrable. The systems under study possessed time-reversal symmetry and so their states should change sign round degeneracies. Seeking to display some degeneracies and their sign changes, M. Wilkinson and I [44] made a detailed investigation of the spectra of vibrating triangles as a function of angles (two parameters).

After a seminar reporting this work in the spring of 1983 at the Georgia Institute of Technology, R. Fox asked me, “what happens to the sign change if a magnetic field is switched on?”, and this question led directly to the discovery of the phase and its 2-form several weeks later. Only when the work was written in first draft was I made aware (by E. Heller) of the papers by Mead and Truhlar. In August 1983, after my paper [7] had been submitted for publication, I described the phase to B.Simon, who instantly saw its relationship to Hermitian line bundles and Chern classes. His paper [10] directed many people towards this subject, thereby provoking the considerable activity of which this book is a partial record. But thanks to a referee’s delay and an accident of astronomy, his paper appeared in 1983, mine in 1984.

Acknowledgments

I thank John Hannay for many inspirations in discussions of phase matters over several years. No military agency supported this research.

References

- [1] Landau, L.D., Lifshitz, E.M., and Pitaevskii, L.P. 1984 *Electrodynamics of Continuous Media, 2nd ed.* (vol. 8 of Course of Theoretical Physics). Oxford: Pergamon Press.
- [2] Born, M., and Wolf, E. 1956 *Principles of Optics*. London: Pergamon Press.
- [3] Chiao, R.Y. and Wu, Y.S. 1986 *Phys. Rev. Lett.* **57**, 933–936. Tomita, A. and Chiao, R.Y. 1986 *Phys. Rev. Lett.* **57**, 937–940.
- [4] Segert, J. 1987 *Phys. Rev.* **A36**, 10–16.
- [5] Haldane, F.D.M. 1896 *Optics Letters* **11**, 730-732.
- [6] Berry, M.V. 1987 *Nature* **326**, 277-278.
- [7] Berry, M.V. 1984 *Proc.Roy.Soc.London* **A392**, 45-57.
- [8] Aharonov, Y. and Anandan, J. 1987 *Phys. Rev. Lett.* **58**, 1593–1596.
- [9] Born, M. and Fock, V. 1928 *Z.Phys.* **51**, 165–169.
- [10] Simon, B. 1983 *Phys. Rev. Lett.* **51**, 2167–2170.
- [11] Bitter, T. and Dubbers, D. 1987 *Phys. Rev. Lett.* **59**, 251-254.
- [12] Tycko, R. 1987 *Phys. Rev. Lett.* **58**, 2281–2284.
- [13] Hannay, J.H. 1985 *J.Phys.* **A18**, 221-230.
- [14] Arnold, V.I. 1978 *Mathematical Methods of Classical Mechanics*. New York: Springer.
- [15] Berry, M.V. 1986 “Adiabatic Phase Shifts for Neutrons and Photons,” in *Fundamental Aspects of Quantum Theory*, eds. V.Gorini and A.Frigerio, NATO ASI series vol.144, 267-278. New York: Plenum.
- [16] Cina, J. 1986 *Chem.Phys.Lett.* **132**, 393-95.
- [17] Littlejohn, R.G. 1984 *Contemp.Math.* **28**, 151.
- [18] Berry, M.V. 1985 *J.Phys.* **A18**, 15-27.
- [19] Berry, M.V. 1983 *Semiclassical Mechanics of Regular and Irregular Motion*, in *Chaotic Behavior of Deterministic Systems*. Les Houches Lecture Series XXXVI, eds G.Iooss, R.H.G.Helleman and R.Stora. Amsterdam: North-Holland. pp.171-271.
- [20] Pancharatnam, S. 1956 *Proc.Ind.Acad.Sci.* **A44**, 247-262. Pancharatnam, S. 1975 *Collected Works of S Pancharatnam*, Oxford: University Press.

- [21] Berry, M.V. 1987 *J.Mod.Optics* **34**, 1401-1407.
- [22] Poston, T. and Stewart, I.N. 1978 *Catastrophe Theory and its Applications*. London: Pitman.
- [23] Berry, M.V. and Upstill, C. 1980 *Prog.Optics* **18**, 257-346.
- [24] Mead, C.A. 1980 *Chem.Phys.(Netherlands)* **49**, 23-32, 33-38.
- [25] Moody, J., Shapere, A., and Wilczek, F. 1986 *Phys.Rev.Lett.* **56**, 893-896.
- [26] Jackiw, R. 1988 *Comm.At.Mol.Phys.* **20**, 71.
- [27] Zygelman, B. 1987 *Phys.Lett.A.* **125**, 476-481.
- [28] Anandan, J. and Aharonov, Y. 1988 *Phys. Rev. Lett.*, in press.
- [29] Berry, M.V. 1987 *Proc.Roy.Soc.* **A414**, 31-46.
- [30] Garrison, J.C. 1986 Preprint UCRL-94267 from Lawrence Livermore Laboratory.
- [31] Ehrenfest, P. 1916 *Ann.d.Physik* **51**, 327-352.
- [32] Berry, M.V. and Mount, K.E. 1972 *Reps.Prog.Phys.* **35**, 315-397.
- [33] Page, D.H. 1987 *Phys.Rev.* **A36**, 3479-3481.
- [34] Porter, C.E. 1965 *Statistical Theories of Spectra: Fluctuations*. New York: Academic Press.
- [35] Uhlenbeck, K. 1976 *Am.J.Math.* **98**, 1059-1078.
- [36] Herzberg, G. and Longuet-Higgins, H.C. 1963 *Disc.Far.Soc.* **35**, 77-82.
- [37] Darboux, G. 1986 *Leçons sur la Théorie Générale des Surfaces*, vol.4, Paris: Gauthier-Villars, note VII.
- [38] Porteous, I.R. 1971 *J.Diff.G geom.* **5**, 543-564.
- [39] Berry, M.V. and Hannay, J.H. 1977 *J.Phys.* **A10**, 1809-1821.
- [40] Ramaseshan, S. and Nityananda, R. 1986 *Current Science (India)* **55**, 1225-26.
- [41] Bouchiat, C. and Gibbons, G.W. 1988 *J.Phys. France* **49**, 187-199.
- [42] Mead, C.A. and Truhlar, D.G. 1979 *J.Chem.Phys.* **70**, 2284-2296.
- [43] Berry, M. 1981 *Ann.Phys.(N.Y.)* **131**, 163-216.
- [44] Berry, M.V. and Wilkinson, M. 1984 *Proc.Roy.Soc.Lond.* **A392**, 15-43.
- [45] Wilkinson, M. 1984 *J.Phys.* **A17**, 3459-3476.
- [46] Wilczek, F. and Zee, A. 1984 *Phys. Rev. Lett.* **52**, 2111-2114.
- [47] Dingle, R.B. 1973 *Asymptotic Expansions: Their Derivation and Interpretation*, New York: Academic Press.
- [48] Frank, F.C. 1958 *Faraday Soc.Disc.* **25**, 19-28.
- [49] Bender, C.M. and Papanicolaou, N. 1988 *J.Phys.France* **49**, 561-566.
- [50] Provost, J.P. and Vallee, G. 1980 *Comm.Math.Phys.* **76**, 289-301.

THREE ELABORATIONS ON BERRY'S CONNECTION, CURVATURE AND PHASE*†

R. JACKIW

*Center for Theoretical Physics, Laboratory for Nuclear Science and Department of Physics,
 Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 U.S.A.*

Received 23 October 1987

We discuss how symmetries and conservation laws are affected when Berry's phase occurs in a quantum system: symmetry transformations of coordinates have to be supplemented by gauge transformations of Berry's connection, and consequently constants of motion acquire terms beyond the familiar kinematical ones. We show how symmetries of a problem determine Berry's connection, curvature and, once a specific path is chosen, the phase as well. Moreover, higher order corrections are also fixed. We demonstrate that in some instances Berry's curvature and phase can be removed by a globally well-defined, time-dependent canonical transformation. Finally, we describe how field theoretic anomalies may be viewed as manifestations of Berry's phase.

We frequently analyze a physical system that is naturally and conveniently divided into two parts. We deal first with the motion of one set of variables, keeping the others fixed but arbitrary, and then complete the analysis of the whole by allowing variation of the previously fixed coordinates. The initially fixed variables we shall call *slow*, those whose motion is analyzed first are the *fast* variables—the terminology reflects the fact that molecular physicists and quantum chemists have used this decomposition for a long time in their Born–Oppenheimer studies of molecules. After resolving the dynamics of the fast variables, one is left with an effective action governing the slow variables. It has been established that this effective dynamics frequently involves an external vector potential \mathbf{A} which is induced by the fast variables. Moreover, the effective equation of motion satisfied by the slow variables involves only the curl of the vector potential—a magnetic-like field \mathbf{B} . Consequently, there is a gauge invariance in the description and only the gauge invariant portion of \mathbf{A} leads to physical effects. This was first seen in Born–Oppenheimer studies of molecules,¹ and analogous effects were also found in quantum field theory.² With Berry's beautiful analysis,³ we appreciate the full quantum mechanical generality of the phenomenon. The induced vector potential \mathbf{A} is now called *Berry's connection*, the induced magnetic-like field \mathbf{B} is *Berry's curvature*, while a line integral of the connection is *Berry's phase*.

* This work is supported in part by funds provided by the U.S. Department of Energy (D.O.E.) under contract number DE-AC02-76ER03069.

† Lecture delivered at “Non-Integrable Phase in Dynamical Systems,” Theoretical Physics Institute, University of Minnesota, Minneapolis, Minnesota, October 1987.

Here, I shall speak on three topics: how symmetries and conservation laws are affected by Berry's connection; how higher order effects, beyond Berry's, may be determined with the help of symmetries; and finally, how Berry's connection manifests itself in quantum field theory through the anomaly phenomenon.

Let me begin by setting some notation. The complete Hamiltonian for the fast (\mathbf{p}, \mathbf{r}) and slow (\mathbf{P}, \mathbf{R}) variables is

$$H = \frac{\mathbf{P}^2}{2M} + \frac{\mathbf{p}^2}{2\mu} + V(\mathbf{R}, \mathbf{r}). \quad (1)$$

The sub-Hamiltonian governing the fast variables depends parametrically on the slow coordinates

$$h(\mathbf{R}) = \frac{\mathbf{p}^2}{2\mu} + V(\mathbf{R}, \mathbf{r}), \quad (2)$$

and possesses "instantaneous" eigenfunctions $|n; \mathbf{R}\rangle$ and eigenvalues $\varepsilon_n(\mathbf{R})$, which also depend on \mathbf{R}

$$h(\mathbf{R})|n; \mathbf{R}\rangle = \varepsilon_n(\mathbf{R})|n; \mathbf{R}\rangle. \quad (3)$$

These states can give rise to Berry's connection,³ defined by

$$\mathbf{A}_n(\mathbf{R}) = \langle n; \mathbf{R} | i\nabla_{\mathbf{R}} | n; \mathbf{R} \rangle. \quad (4)$$

In the above, I have taken the n^{th} eigenstate to be non-degenerate. Otherwise, one works within the degenerate subspace, and Berry's connection becomes a matrix in this space—an induced non-Abelian Yang-Mills-like potential.⁴

The effective Hamiltonian for slow motion, in the Born–Oppenheimer approximation, which results when off-diagonal matrix elements $\langle n; \mathbf{R} | i\nabla_{\mathbf{R}} | n'; \mathbf{R} \rangle$ are ignored, reads

$$H_{\text{eff}} = \frac{1}{2M}(\mathbf{P} - \mathbf{A}_n(\mathbf{R}))^2 + \varepsilon_n(\mathbf{R}). \quad (5)$$

Correspondingly, the effective Lagrangian is

$$L_{\text{eff}} = \frac{1}{2}M\dot{\mathbf{R}}^2 + \dot{\mathbf{R}} \cdot \mathbf{A}_n(\mathbf{R}) - \varepsilon_n(\mathbf{R}). \quad (6)$$

Equations (5) and (6) show that the fast system induces into the slow dynamics a potential energy $\varepsilon_n(\mathbf{R})$, and a velocity dependent interaction arising from Berry's connection \mathbf{A}_n , which enters dynamics only through the curvature \mathbf{B}_n associated with \mathbf{A}_n .

$$\mathbf{B}_n = \nabla \times \mathbf{A}_n - i\mathbf{A}_n \times \mathbf{A}_n. \quad (7)$$

The last term is present in the non-Abelian matrix case.

Gauge invariance of physical processes means that physical content is not affected by gauge transformations that change \mathbf{A} according to

$$\mathbf{A}(\mathbf{R}) \rightarrow g^{-1}(\mathbf{R})\mathbf{A}(\mathbf{R})g(\mathbf{R}) - ig^{-1}(\mathbf{R})\nabla g(\mathbf{R}). \quad (8)$$

[Henceforth we suppress the label “ n ” on the induced quantities.] Transformation (8) arises when the instantaneous eigenfunctions are redefined by R -dependent phase factors, which are not fixed by the instantaneous eigenvalue equation (3). Formula (8) is presented for the general case, where the connection \mathbf{A} is matrix valued, belonging to the Lie algebra of some group, and g is the matrix representation for a group element. For example, \mathbf{A} could be a Hermitian 2×2 matrix, and then g is a unitary 2×2 matrix belonging to $U(2)$. In the simplest, Abelian case, \mathbf{A} is a 1×1 matrix—a function—everything commutes and g belonging to $U(1)$ is given by $e^{i\theta}$, so that (8) reduces to the electromagnetic gauge transformation

$$\mathbf{A} \rightarrow \mathbf{A} + \nabla\theta. \quad (9)$$

The gauge transformation (8) may also be presented in infinitesimal form. Writing

$$g = I + i\Theta + \dots \quad (10)$$

where Θ is a Hermitian matrix in the Lie algebra, we deduce from (8) that the infinitesimal transformation is

$$\delta\mathbf{A} = \nabla\Theta - i[\mathbf{A}, \Theta]. \quad (11)$$

In the Abelian case, when \mathbf{A} and Θ are functions, the commutator vanishes, and (11) reduces to (9).

We are now at our first topic: Let us suppose that the full Hamiltonian (1) possesses a symmetry and consequently there exist constants of motion commuting with H . We expect that in the effective description, the symmetry is not lost and there are effective constants of motion that commute with H_{eff} . The question is how do the symmetries manifest themselves and how are the constants of motion modified in the presence of Berry's connection.

Let us for definiteness concentrate on rotational symmetry and on the associated angular momentum constant of motion \mathbf{J} . Under rotations the coordinate \mathbf{R} transforms according to $R^i \rightarrow \Lambda^i R^j$ with Λ^i a special orthogonal matrix. We certainly know what it means for the effective potential energy to be invariant against rotations: $\epsilon(\mathbf{R})$ should depend only on the magnitude $R = |\mathbf{R}|$. Analytically, this requirement of rotational invariance is presented as

$$\varepsilon(\Lambda \mathbf{R}) = \varepsilon(\mathbf{R}). \quad (12)$$

For a vector quantity, like Berry's connection \mathbf{A} , the analogous requirement recognizes the vectorial nature of \mathbf{A}

$$\mathbf{A}(\Lambda \mathbf{R}) = \Lambda \mathbf{A}(\mathbf{R}). \quad (13)$$

Equivalently,

$$\Lambda \mathbf{A}(\Lambda^{-1} \mathbf{R}) = \mathbf{A}(\mathbf{R}). \quad (14)$$

Vector functions satisfying (14) have the form $\mathbf{A}(\mathbf{R}) = \mathbf{R}A(R)$, where A is a scalar function of $R = |\mathbf{R}|$. However, such connections are not physically interesting because the associated curvature (7)—the magnetic-like field—vanishes.

For a nontrivial realization of rotational invariance, we need to recognize that symmetries should be looked for in the physical content of the formalism, and we must remember that the connection has unphysical components, because a gauge transformation can always be performed without affecting physical content.

The rotational invariance requirement for a vector connection can be weakened from (14), while still retaining the full force of rotational symmetry for physical processes. Rather than demanding that $\Lambda \mathbf{A}(\Lambda^{-1} \mathbf{R})$ equal $\mathbf{A}(\mathbf{R})$, as in (14), we allow a gauge transformation to intervene.

$$\Lambda \mathbf{A}(\Lambda^{-1} \mathbf{R}) = g^{-1}(\mathbf{R}) \mathbf{A}(\mathbf{R}) g(\mathbf{R}) - i g^{-1}(\mathbf{R}) \nabla g(\mathbf{R}). \quad (15)$$

For most purposes, it is sufficient to consider the infinitesimal version of (15). When we define

$$\Lambda^i = \delta^{ij} - \epsilon^{ijk} n^k + \dots \quad (16a)$$

so that an infinitesimal rotation of \mathbf{R} is

$$\delta \mathbf{R} = \mathbf{n} \times \mathbf{R}, \quad (16b)$$

then (11) and (15) combine to

$$\mathbf{n} \times \mathbf{A} - (\mathbf{n} \times \mathbf{R} \cdot \nabla) \mathbf{A} = \nabla \Theta - i[\mathbf{A}, \Theta]. \quad (17)$$

Equation (17) may also be presented as a condition on the curvature \mathbf{B} , (7). An algebraic rearrangement of (17), which uses various vector identities, casts (17) in the form

$$(\mathbf{n} \times \mathbf{R}) \times \mathbf{B} = \nabla W - i[\mathbf{A}, W], \quad (18a)$$

$$W = \Theta + \mathbf{n} \times \mathbf{R} \cdot \mathbf{A}. \quad (18b)$$

To summarize: the effective Born–Oppenheimer Hamiltonian (5) and the effective Born–Oppenheimer Lagrangian (6) are rotationally invariant, provided, in addition to (12)—which expresses rotational symmetry of the instantaneous energy—the Berry connection \mathbf{A} satisfies (15) or (17), equivalently the Berry curvature satisfies (18)—this assures rotational symmetry up to a gauge transformation of the connection. Moreover, one can prove with general quantum mechanical reasoning that when the fast Hamiltonian $h(\mathbf{R})$ is invariant against rotations of the fast variables, supplemented by a rotation of the slow background variable \mathbf{R} , then the instantaneous eigenstates $|n; \mathbf{R}\rangle$ transform under rotations precisely in such a way that the connection defined in (4) satisfies (15).⁵ Therefore, as expected, rotational symmetry is not lost in the effective Hamiltonian, but is realized in a nontrivial fashion, when a gauge connection is present.

In the known models with nonvanishing Berry's connection, one can verify (17) or (18). Consider the example,³

$$h(\mathbf{R}) = \mathbf{S} \cdot \mathbf{R}(t) \quad (19)$$

where \mathbf{S} , playing the role of the fast variables (\mathbf{p}, \mathbf{r}) , is an angular momentum operator that obeys the SO(3) Lie algebra

$$[S_i, S_j] = i\epsilon_{ijk} S^k. \quad (20)$$

h is invariant against simultaneous rotations of \mathbf{S} and of \mathbf{R} , the slow background. The instantaneous eigenvalues are mR , with m ranging in unit steps from $-S$ to S . The model gives rise to Berry's curvature, which is proportional to a magnetic monopole field³

$$\mathbf{B} = -m \frac{\hat{\mathbf{R}}}{R^2}. \quad (21)$$

It is straightforward to verify that (18) is satisfied with

$$W = m\hat{\mathbf{R}} \cdot \mathbf{n}. \quad (22)$$

\mathbf{B} in Eq. (21) is an Abelian curvature. In a non-Abelian U(2) example, relevant to diatoms with λ degeneracy, Berry's connection is a 2×2 matrix, which depends on a parameter κ ⁶

$$\mathbf{A}(\mathbf{R}) = \frac{1 + \kappa}{2} \frac{\hat{\mathbf{R}} \times \boldsymbol{\sigma}}{R}. \quad (23)$$

The $\boldsymbol{\sigma}$ matrices are Pauli matrices. Rotational invariance holds: \mathbf{A} satisfies (17) and Θ is given by⁷

$$\Theta = \frac{1}{2} \boldsymbol{\sigma} \cdot \mathbf{n}. \quad (24)$$

In fact one may reverse the reasoning and determine the connection or curvature from the requirements of rotational symmetry. One views (18) as an equation for \mathbf{B} and finds the most general solution in terms of W . The point is that integrability requirements on (18) sharply limit the allowed solutions. Thus the magnetic monopole in (21) is the *unique* Abelian rotationally symmetric curvature, while the configuration (23) is one of the few non-Abelian 2×2 matrix solutions. In this way knowledge of the symmetry limits the curvature, apart from the gauge freedom.

A further remark: because invariance of H_{eff} is achieved when a rotation is supplemented by a gauge transformation, constants of motion commuting with H_{eff} are modified.⁸ Recall that a constant of motion is also the generator of the infinitesimal transformation. In our case the infinitesimal rotation of coordinates must be supplemented by an infinitesimal gauge transformation of the connection. Consequently, the conserved angular momentum \mathbf{J} arising from rotational symmetry in the $\hat{\mathbf{n}}$ direction is the conventional $\mathbf{n} \cdot (\mathbf{R} \times \mathbf{P})$ supplemented by Θ , the gauge transformation generator. Equivalently,

$$\mathbf{J} = \mathbf{R} \times M\dot{\mathbf{R}} + \frac{\partial W}{\partial \mathbf{n}} \quad (25)$$

where we have used (18b) and the fact that $\mathbf{P} = M\dot{\mathbf{R}} + \mathbf{A}$. [Of course, Θ and W are linear in \mathbf{n} .] The extra term in (25) puts into evidence the physical significance of the gauge transformation that accompanies the rotational coordinate change: $\frac{\partial W}{\partial \mathbf{n}}$ is the angular momentum stored in the gauge field, that summarizes the effect of the fast variables on the slow ones. This is the reason why in the presence of a magnetic monopole [or of the Berry curvature (21)] the angular momentum contains in addition to the kinematical term the further $m\dot{\mathbf{R}}$ [see (22)]. Also this is why the U(2) non-Abelian connection (23) requires that the kinematical angular momentum be supplemented by $\sigma/2$ [see (24)].

We have discussed in detail how rotational invariance and its associated conserved quantity/generator are affected by Berry's connection. However the same ideas can be applied to invariance under arbitrary coordinate transformations. When the complete Hamiltonian is invariant against some coordinate transformation, the effective Hamiltonian for slow motion retains this property, but the background gauge field induced by the fast variables is invariant only up to a gauge transformation,⁵ whose infinitesimal generator supplements the relevant constant of motion.⁸ The problem of determining all connections/curvatures that are invariant, up to a gauge transformation, against an arbitrary coordinate transformation has been solved by mathematicians and physicists.⁹

There is one more model worth mentioning: the generalized one-dimensional harmonic oscillator, with time-varying parameters

$$h = \frac{1}{2}(\alpha(t)p^2 + \beta(t)(pq + qp) + \gamma(t)q^2), \quad (26)$$

$\alpha > 0.$

Berry's connection is nonvanishing.¹⁰ The three operators p^2 , $pq + qp$ and q^2 close upon commutation on the Lie algebra of $\text{SO}(2, 1)$, which is like the angular momentum $\text{SO}(3)$ algebra, except some crucial signs are reversed. This may be presented by redefining the fast variables as

$$T_1 = \frac{1}{4}(q^2 - p^2), \quad T_2 = \frac{1}{4}(pq + qp), \quad T_3 = \frac{1}{4}(q^2 + p^2) \quad (27)$$

and the parameters

$$R^1 = \gamma - \alpha, \quad R^2 = 2\beta, \quad R^3 = \gamma + \alpha. \quad (28)$$

The Hamiltonian (26) takes the form

$$h = T_i R^i(t), \quad (29)$$

with

$$[T_i, T_j] = i\epsilon_{ijk} T^k. \quad (30)$$

Formulas (29) and (30) are analogous to (19) and (20) of the $\text{SO}(3)$ case, except that here a metric g_{ij} intervenes in the ijk group indices, and introduces the crucial sign differences from $\text{SO}(3)$

$$g_{ij} = \text{diag}(-1, -1, 1). \quad (31)$$

The general discussion applies, and the curvature is immediately predicted to be the unique $\text{SO}(2, 1)$ covariant,⁵

$$B^i \propto \frac{R^i}{(R^j R_j)^{3/2}} \quad (32)$$

where $R^j R_j = (R^3)^2 - (R^1)^2 - (R^2)^2 = 4(\alpha\gamma - \beta^2)$ is assumed to be positive. Equation (32) coincides with the result of the explicit calculation based on (25).¹⁰

In spite of the formal similarity to the $\text{SO}(3)$ problem, there is a crucial difference which has not been remarked upon in the literature. The curvature may be removed, not of course, by a gauge transformation, but by a globally well-defined canonical transformation, or what is equivalent by dropping a total time derivative from the Lagrangian that corresponds to the Hamiltonian (26)

$$L \equiv pq - h = \frac{1}{2\alpha} \dot{q}^2 - \frac{1}{2} \left(\gamma - \frac{\beta^2}{\alpha} - \frac{d}{dt} \frac{\beta}{\alpha} \right) q^2 - \frac{d}{dt} \frac{1}{2} \left(\frac{\beta}{\alpha} q^2 \right). \quad (33)$$

When the last term is dropped, as it can be without affecting dynamics, one is left with an equivalent theory, described by the Hamiltonian

$$\bar{h} = \frac{\alpha}{2} p^2 + \frac{1}{2} \left(\gamma - \frac{\beta^2}{\alpha} - \frac{d}{dt} \frac{\beta}{\alpha} \right) q^2, \quad (34)$$

which does not lead to a Berry connection, but is canonically equivalent to h , with the time-dependent canonical transformation $p \rightarrow p + \frac{\beta}{\alpha} q$, $q \rightarrow q$. Notice $\frac{\beta}{\alpha}$ is non-singular, because α is assumed never to vanish. Hence this canonical transformation is globally defined on the parameter space. On the other hand, an analogous transformation for the SO(3) example,³ which would rotate $\mathbf{S} \cdot \mathbf{R}(t)$ into a fixed direction, cannot be globally defined.

The reason for this difference between the SO(3) and SO(2, 1) examples derives from the fact that the parameter space [at fixed R] of the SO(3) model is the surface of a sphere, while that of the SO(2, 1) model is one sheet of the hyperboloid $R^i R_i = \text{constant}$. Unlike the former, the latter is topologically trivial, and homotopic to the Euclidean plane.¹¹

Symmetry considerations may also be used to calculate quickly and efficiently higher order corrections to Berry's phase, and this brings us to the second topic in my lecture. First we need to define what we wish to calculate in higher order. I have already indicated in Eq. (6) that the instantaneous energy eigenvalue and Berry's connection are two contributions to the effective Lagrangian induced by the fast variables onto the slow ones. Let us therefore consider the complete effective action I_{eff} induced by the fast variables. This quantity may be defined as follows. Consider the time-dependent Schrödinger equation with the time-dependent Hamiltonian $h(\mathbf{R}(t))$

$$i\partial_t |\psi; t\rangle = h(\mathbf{R}(t)) |\psi; t\rangle. \quad (35)$$

We take the two [in, out] solutions of (35) $|\psi^{(\pm)}; t\rangle$ that satisfy the initial and final conditions, respectively

$$\lim_{t \rightarrow t_i} |\psi^{(+)}; t\rangle = |n; \mathbf{R}(t_i)\rangle, \quad (36a)$$

$$\lim_{t \rightarrow t_f} |\psi^{(-)}; t\rangle = |n; \mathbf{R}(t_f)\rangle. \quad (36b)$$

For simplicity we shall assume $\mathbf{R}(t_i) = \mathbf{R}(t_f)$, but this can be relaxed. The effective action is defined from the in-out matrix element,

$$I_{\text{eff}} = -i \ln \langle \psi^{(-)}; t_i | \psi^{(+)}; t_f \rangle. \quad (37)$$

[Alternatively, I_{eff} may be given by a Feynman path integral over the fast variables.] I_{eff} is a functional of $\mathbf{R}(t)$, but it is independent of t because $|\psi^\pm; t\rangle$ satisfy the Schrödinger equation. Hence the overlap in (37) may be evaluated variously at $t = t_i$, where $I_{\text{eff}} = -i \ln \langle \psi^{(-)}; t_i | n; \mathbf{R}(t_i) \rangle$ or at $t = t_f$, where $I_{\text{eff}} = -i \ln \langle n; \mathbf{R}(t_f) | \psi^{(+)}; t_f \rangle$.

One may expand I_{eff} in a series of terms with increasing time derivatives of \mathbf{R} . The zeroth order term has no time derivatives; it is what would survive of I_{eff} when $\mathbf{R}(t)$ is time-independent. The first order term has a single time derivative; the second order term involves two time derivatives, etc. We now see that the instantaneous energy eigenvalue is the (negative) zeroth order term; the Berry phase gives the first order term, linear in time derivatives; higher orders are to be determined—compare (6).

$$\begin{aligned} I_{\text{eff}} = & - \int dt \varepsilon(\mathbf{R}(t)) + \int dt \dot{\mathbf{R}}(t) \cdot \mathbf{A}(\mathbf{R}(t)) \\ & + \frac{1}{2} \int dt \dot{\mathbf{R}}^i(t) M^{ij}(\mathbf{R}(t)) \dot{\mathbf{R}}^j(t) + \dots \end{aligned} \quad (38)$$

The coefficient tensors A^i and M^{ij} must be transverse to R^i . This is so because a purely radial time dependence, $\mathbf{R}(t) = \hat{\mathbf{R}}R(t)$ with time-independent $\hat{\mathbf{R}}$, gives rise only to the first term in (38). Equation (38) is an asymptotic series, it produces a real I_{eff} but misses imaginary parts, which describe decay. [Note that a second order term of the form $\int dt \dot{\mathbf{R}}(t) \cdot \mathbf{N}(\mathbf{R}(t))$ is equivalent, by an integration by parts, to the last term in (38); to justify dropping endpoint contributions, the motion must be periodic, alternatively take $t_{i,f} = \mp\infty$, where everything vanishes.]

Let me show how all higher order terms in (38) can be computed from the knowledge of the zeroth order, instantaneous eigenvalue. The discussion is confined to the SO(3) example. Observe that the generator of rotations on the fast variables \mathbf{S} , is obtained from $h(\mathbf{R})$ by differentiating it with respect to \mathbf{R} ,

$$\mathbf{S} = \frac{\partial h(\mathbf{R})}{\partial \mathbf{R}}. \quad (39)$$

This operator satisfies,

$$i[h(\mathbf{R}), \mathbf{S}] = \mathbf{R} \times \mathbf{S}, \quad (40a)$$

which is the Heisenberg picture equation for $\mathbf{S}(t)$,

$$\frac{d}{dt} \mathbf{S}(t) - \mathbf{R}(t) \times \mathbf{S}(t) = 0. \quad (40b)$$

The above non-conservation equation for $\mathbf{S}(t)$ reflects the fact that rotations on the fast variables \mathbf{S} , are not symmetry operations, unless supplemented by a rotation of the external parameters.

Next I define the in-out matrix element of \mathbf{S} ,

$$\mathbf{J}(t) = \langle \psi^{(-)}; t | \mathbf{S} | \psi^{(+)}; t \rangle / \langle \psi^{(-)}; t | \psi^{(+)}; t \rangle \quad (41)$$

which by virtue of (40a) also satisfies (40b). Finally observe that $\mathbf{J}(t)$ is given by a functional derivative with respect to $\mathbf{R}(t)$ of I_{eff} [compare with (39)],

$$\mathbf{J}(t) = -\frac{\delta I_{\text{eff}}}{\delta \mathbf{R}(t)}. \quad (42)$$

Thus combining (40) with (42) we arrive at a condition on I_{eff}

$$\frac{d}{dt} \frac{\delta I_{\text{eff}}}{\delta \mathbf{R}(t)} - \mathbf{R} \times \frac{\delta I_{\text{eff}}}{\delta \mathbf{R}(t)} = 0. \quad (43)$$

Equation (43) may be applied iteratively to (38), and then equal orders of time derivatives of $\mathbf{R}(t)$ are equated. Varying I_{eff} gives

$$\begin{aligned} \mathbf{J}(t) &= m\hat{\mathbf{R}}(t) - \dot{\mathbf{R}}(t) \times \mathbf{B}(\mathbf{R}(t)) \\ &- \frac{\delta}{\delta \mathbf{R}(t)} \frac{1}{2} \int d\tau \dot{\mathbf{R}}^i(\tau) M^{ij}(\mathbf{R}(\tau)) \dot{\mathbf{R}}^j(\tau) + \dots \end{aligned} \quad (44)$$

and (43) requires

$$m \frac{d}{dt} \hat{\mathbf{R}} = -\mathbf{R} \times (\dot{\mathbf{R}} \times \mathbf{B}), \quad (45a)$$

$$\frac{d}{dt} (\dot{\mathbf{R}} \times \mathbf{B}(\mathbf{R})) = \mathbf{R} \times \frac{\delta}{\delta \mathbf{R}} \frac{1}{2} \int d\tau \dot{\mathbf{R}}^i M^{ij} \dot{\mathbf{R}}^j, \quad (45b)$$

etc. Equation (45a) determines \mathbf{B} : the known formula (21) is regained; Eq. (45b) gives a new result: the second order term in the derivative expansion, which is the first correction to Berry's phase. One finds

$$M^{ij} = \frac{m}{R^3} (\delta^{ij} - \hat{\mathbf{R}}^i \hat{\mathbf{R}}^j). \quad (46)$$

Clearly higher order terms may be similarly computed. The result (46) has been verified by a direct evaluation of the in, out matrix element.¹² Of course a specific value for I_{eff} is obtained only when a specific path $\mathbf{R}(t)$ is chosen and the time integral is performed.

My third and last topic concerns the role that Berry's phase plays in modern quantum field theory, where it gives another point of view on the anomaly phenomenon. This subject will be discussed in greater detail in another lecture.¹³ Here I give an introductory description, because I believe that this peculiar feature of second quantized field theories is probably unfamiliar to many of you.

I begin by reminding that the quantum mechanical revolution has not erased our reliance on the earlier classical physics. Indeed when proposing a theory, we begin with classical concepts and construct models according to the rules of classical, pre-quantum physics. We know, however, such classical reasoning is not in complete accord with quantum reality. Therefore, the classical model is reanalyzed by the rules of quantum physics, which comprise the true laws of nature, i.e., the model is *quantized*. For a long time it was believed that symmetries of a theory are not affected by the transition from classical to quantum rules. However, more recently we have learned that this is not so. In a quantized theory, some symmetries of classical physics may disappear because symmetry violating processes, which are not seen classically, can occur when the analysis is conducted with quantum effects taken into consideration. Such tenuous symmetries are said to be *anomalously broken*. Although present classically, they are absent from the quantum version of the theory, unless the model is carefully arranged to avoid this effect.

Anomalously or quantum mechanically broken symmetries play a crucial role in our present-day theories of elementary particles. In some instances they save the models from possessing too much symmetry, which would not be in accord with experiment. In other instances, the desire to preserve a symmetry in the quantum theory places strong constraints on model building and gives experimentally verified predictions. For example, the equality in the number of quarks and leptons is understood in these terms. Also, the present-day excitement about strings derives from the fact that only very few string models can be adjusted to avoid quantum mechanical, anomalous breaking of those symmetries that make string theory free of the infinities plaguing conventional field theories. Thus the number of consistent string models appears very limited, and a limitation of theoretical possibilities is what every model builder looks for. Anomalous symmetries are also beginning to play a role in other branches of physics, like condensed matter.

For a specific example of this phenomenon, consider massless fermions moving in an electromagnetic field described by electromagnetic potentials. Since massless fermions possess a well-defined helicity, we shall consider fermions with only one helicity. Such systems are an ingredient in theories of quarks and leptons. Moreover, they also arise in condensed matter physics, not because one is dealing with massless, single-helicity particles, but because a well-formulated approximation to some many-body Hamiltonians can result in a first order matrix equation which is identical to the equation for single-helicity massless fermions, i.e. a massless Dirac equation for a spinor ψ .

As a first quantized theory, the system is gauge covariant, in that a gauge transformation on the electromagnetic potential can be compensated by a change of the wavefunction, ψ . Moreover, the norm of the wavefunction is time-independent: $N = \int \psi^\dagger \psi, \frac{d}{dt} N = 0$. So far there are no surprises.

To construct a quantum field theory from the above, the model is second quantized: the wavefunction ψ is promoted to an operator Ψ and the state space for the theory is a many-particle Fock space. Moreover, the ground state of the second quantized

field theory has to be the filled *Dirac sea* so that the negative energy solutions of the first quantized Dirac equation are eliminated. Of course all states are functionals of the background electromagnetic potential in which the fermions move.

One expects that the second quantized theory also possesses gauge invariance, and that as a consequence of gauge invariance the total charge to which the first quantized norm is promoted, $Q = \int \Psi^\dagger \Psi$, is conserved. In fact, this is not true: the charge is not conserved; rather one finds

$$\frac{d}{dt} Q \propto \int \mathbf{E} \cdot \mathbf{H} \quad (47)$$

where \mathbf{E} and \mathbf{H} are the background electric and magnetic fields, and correspondingly gauge invariance is lost.¹⁵

There are many ways of arriving at the result where the gauge symmetry in this model is anomalously broken. In one physically transparent argument, it is established that the process of filling the negative energy sea to define the field theoretic ground state necessarily violates gauge invariance.¹⁴

Berry's ideas provide another viewpoint. The fermion field operators Ψ are thought of as fast variables, the analogs of \mathbf{p} and \mathbf{r} , or of \mathbf{S} in the example (19). The background electromagnetic potential is viewed as an external parameter; it is the analog of \mathbf{R} . The Fock-space state vectors are functionals of the background potential; they are analogs of $|n; \mathbf{R}\rangle$. An analysis shows that when the background electromagnetic potentials are gauge transformed—this can be an adiabatic change—the states acquire a phase variation and in this way lose electromagnetic gauge invariance. Thus the anomaly phenomenon is a manifestation in quantum field theory of Berry's phase, $\int d\mathbf{R} \cdot \mathbf{A}(\mathbf{R})$.¹⁶

Symmetry in quantum theory can also be seen through the realization of the symmetry algebra in the canonical commutation relations. Correspondingly, when the symmetry is anomalously or quantum mechanically broken, the algebra acquires a dynamical modification. As is seen from (6), the Berry connection induces velocity-dependent interactions, which modify the relation between canonical momentum and velocity. Consequently, the commutator of velocity components acquires a quantum mechanical correction.

$$[M\dot{\mathbf{R}}^i, M\mathbf{R}^j] = i\epsilon^{ijk}B^k \quad (48)$$

Anomalous commutators—an important chapter in the anomaly story¹⁴—are also connected to the Berry curvature.

References

1. H. Longuet-Higgins, U. Opik, M. Pryce and R. Sack, *Proc. Roy. Soc. A***224** (1958) 1; H. Longuet-Higgins, *Adv. Spectrosc.* **2** (1961) 429, *Proc. R. Soc. A***344** (1975) 147; M. Child and H. Longuet-Higgins, *Phil. Trans. R. Soc.* **254** (1961) 259; G. Herzberg and H. Longuet-Higgins, *Disc. Faraday Soc.* **35** (1963) 77; M. O'Brien, *Proc. R. Soc. A***281** (1984) 323; A. Stone, *Proc. R. Soc. A***351** (1976) 141; C. Mead, *J. Chem. Phys.* **70** (1979) 2276, **72** (1980)

Three Elaborations on Berry's Connection, Curvature and Phase 297

- 3839, *Chem. Phys.* **49** (1980) 23, 33; C. Mead and D. Truhlar, *J. Chem. Phys.* **70** (1979) 2284, (E) **78** (1983) 6344.
2. R. Jackiw, in *E. Fradkin Festshrift*, I. Batalin, C. Isham and G. Vilkovisky, eds. (Adam Hilger, Bristol, 1987). See also M. Asorey and D. Mitter, *Phys. Lett.* **153B** (1985) 147; Y.-S. Wu and A. Zee, *Phys. Lett.* **B258** (1985) 157.
 3. M. V. Berry, *Proc. R. Soc. A* **392** (1984) 45, and lecture delivered at this conference.
 4. F. Wilczek and A. Zee, *Phys. Rev. Lett.* **52** (1984) 2111.
 5. L. Vinet, University of Montreal preprint, August (1987).
 6. J. Moody, A. Shapere and F. Wilczek, *Phys. Rev. Lett.* **56** (1986) 893. Formula (23) is presented in a gauge which differs from the one used by Moody et al.
 7. R. Jackiw, *Phys. Rev. Lett.* **56** (1986) 2779.
 8. R. Jackiw and N. Manton, *Ann. Phys. (NY)* **127** (1980) 257.
 9. H. Wang, *Nagoya Math. J.* **13** (1958) 1; J. Harnad, S. Shnider and L. Vinet, *J. Math. Phys.* **21** (1980) 2719; R. Forgács and N. Manton, *Comm. Math. Phys.* **72** (1980) 15. For a review and more details see R. Jackiw, *Acta Phys. Austriaca, Suppl.* **XXII** (1980) 383.
 10. M. V. Berry, *J. Phys. A* **18** (1985) 15; J. Hannay, *J. Phys. A* **18** (1985) 221.
 11. P. Gerbert, MIT preprint CTP #1537, October (1987). See also E. Gozzi and W. Thacker, *Phys. Rev. D* **35** (1987) 2398.
 12. P. Gerbert, Ref. [11]; S. Iida, private communication.
 13. G. Semenoff, lecture delivered at this conference.
 14. For a review, see *Current Algebra and Anomalies*, S. Treiman, R. Jackiw, B. Zumino and E. Witten, eds. (Princeton University Press/World Scientific, Princeton NJ/Singapore, 1985).
 15. D. Gross and R. Jackiw, *Phys. Rev. D* **6** (1972) 477. For details see Ref. [14].
 16. P. Nelson and L. Alvarez-Gaumé, *Comm. Math. Phys.* **99** (1985) 103; H. Sonoda, *Phys. Lett.* **156B** (1985) 220, *Nucl. Phys.* **B266** (1986) 410; A. Niemi and G. Semenoff, *Phys. Rev. Lett.* **55** (1985) 927, **56** (1986) 1019, *Phys. Lett.* **B175** (1986) 439; A. Niemi, G. Semenoff and Y.-S. Wu, *Nucl. Phys.* **B276** (1986) 173. For reviews of this approach to anomalies, see G. Semenoff in *Super Field Theory*, H. Lee, V. Elias, G. Kunstatter, R. Mann and K. Viswanathan, eds. (Plenum, New York, 1987); A. Niemi in *Workshop on Skyrmiions and Anomalies*, M. Jezabek and M. Praszalowicz, eds. (World Scientific, Singapore, 1987).

THE PHYSICAL REVIEW

A journal of experimental and theoretical physics established by E. L. Nichols in 1893

SECOND SERIES, VOL. 115, No. 3

AUGUST 1, 1959

Significance of Electromagnetic Potentials in the Quantum Theory

Y. AHARONOV AND D. BOHM
H. H. Wills Physics Laboratory, University of Bristol, Bristol, England
 (Received May 28, 1959; revised manuscript received June 16, 1959)

In this paper, we discuss some interesting properties of the electromagnetic potentials in the quantum domain. We shall show that, contrary to the conclusions of classical mechanics, there exist effects of potentials on charged particles, even in the region where all the fields (and therefore the forces on the particles) vanish. We shall then discuss possible experiments to test these conclusions; and, finally, we shall suggest further possible developments in the interpretation of the potentials.

1. INTRODUCTION

IN classical electrodynamics, the vector and scalar potentials were first introduced as a convenient mathematical aid for calculating the fields. It is true that in order to obtain a classical canonical formalism, the potentials are needed. Nevertheless, the fundamental equations of motion can always be expressed directly in terms of the fields alone.

In the quantum mechanics, however, the canonical formalism is necessary, and as a result, the potentials cannot be eliminated from the basic equations. Nevertheless, these equations, as well as the physical quantities, are all gauge invariant; so that it may seem that even in quantum mechanics, the potentials themselves have no independent significance.

In this paper, we shall show that the above conclusions are not correct and that a further interpretation of the potentials is needed in the quantum mechanics.

2. POSSIBLE EXPERIMENTS DEMONSTRATING THE ROLE OF POTENTIALS IN THE QUANTUM THEORY

In this section, we shall discuss several possible experiments which demonstrate the significance of potentials in the quantum theory. We shall begin with a simple example.

Suppose we have a charged particle inside a "Faraday cage" connected to an external generator which causes the potential on the cage to alternate in time. This will add to the Hamiltonian of the particle a term $V(x,t)$ which is, for the region inside the cage, a function of time only. In the nonrelativistic limit (and we shall

assume this almost everywhere in the following discussions) we have, for the region inside the cage, $\hat{H} = \hat{H}_0 + V(t)$ where \hat{H}_0 is the Hamiltonian when the generator is not functioning, and $V(t) = e\phi(t)$. If $\psi_0(x,t)$ is a solution of the Hamiltonian \hat{H}_0 , then the solution for \hat{H} will be

$$\psi = \psi_0 e^{-iS/\hbar}, \quad S = \int V(t) dt,$$

which follows from

$$i\hbar \frac{\partial \psi}{\partial t} = \left(i\hbar \frac{\partial \psi_0}{\partial t} + \psi_0 \frac{\partial S}{\partial t} \right) e^{-iS/\hbar} = [\hat{H}_0 + V(t)] \psi = \hat{H} \psi.$$

The new solution differs from the old one just by a phase factor and this corresponds, of course, to no change in any physical result.

Now consider a more complex experiment in which a single coherent electron beam is split into two parts and each part is then allowed to enter a long cylindrical metal tube, as shown in Fig. 1.

After the beams pass through the tubes, they are combined to interfere coherently at E . By means of time-determining electrical "shutters" the beam is chopped into wave packets that are long compared with the wavelength λ , but short compared with the length of the tubes. The potential in each tube is determined by a time delay mechanism in such a way that the potential is zero in region I (until each packet is well inside its tube). The potential then grows as a function of time, but differently in each tube. Finally, it falls back to zero, before the electron comes near the

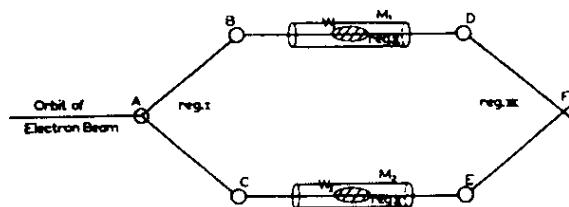


FIG. 1. Schematic experiment to demonstrate interference with time-dependent scalar potential. A, B, C, D, E : suitable devices to separate and divert beams. W_1, W_2 : wave packets. M_1, M_2 : cylindrical metal tubes. F : interference region.

other edge of the tube. Thus the potential is nonzero only while the electrons are well inside the tube (region II). When the electron is in region III, there is again no potential. The purpose of this arrangement is to ensure that the electron is in a time-varying potential without ever being in a field (because the field does not penetrate far from the edges of the tubes, and is nonzero only at times when the electron is far from these edges).

Now let $\psi(x,t) = \psi_1^0(x,t) + \psi_2^0(x,t)$ be the wave function when the potential is absent (ψ_1^0 and ψ_2^0 representing the parts that pass through tubes 1 and 2, respectively). But since V is a function only of t wherever ψ is appreciable, the problem for each tube is essentially the same as that of the Faraday cage. The solution is then

$$\psi = \psi_1^0 e^{-iS_1/\hbar} + \psi_2^0 e^{-iS_2/\hbar},$$

where

$$S_1 = e \int \varphi_1 dt, \quad S_2 = e \int \varphi_2 dt.$$

It is evident that the interference of the two parts at F will depend on the phase difference $(S_1 - S_2)/\hbar$. Thus, there is a physical effect of the potentials even though no force is ever actually exerted on the electron. The effect is evidently essentially quantum-mechanical in nature because it comes in the phenomenon of interference. We are therefore not surprised that it does not appear in classical mechanics.

From relativistic considerations, it is easily seen that the covariance of the above conclusion demands that there should be similar results involving the vector potential, A .

The phase difference, $(S_1 - S_2)/\hbar$, can also be expressed as the integral $(e/\hbar) \oint \varphi dt$ around a closed circuit in space-time, where φ is evaluated at the place of the center of the wave packet. The relativistic generalization of the above integral is

$$\frac{e}{\hbar} \oint \left(\varphi dt - \frac{\mathbf{A}}{c} \cdot d\mathbf{x} \right),$$

where the path of integration now goes over any closed circuit in space-time.

As another special case, let us now consider a path in space only ($t = \text{constant}$). The above argument

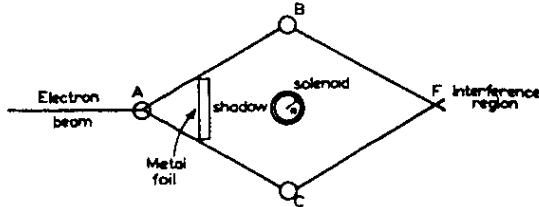


FIG. 2. Schematic experiment to demonstrate interference with time-independent vector potential.

suggests that the associated phase shift of the electron wave function ought to be

$$\Delta S/\hbar = -\frac{e}{ch} \oint \mathbf{A} \cdot d\mathbf{x},$$

where $\oint \mathbf{A} \cdot d\mathbf{x} = \oint \mathbf{H} \cdot d\mathbf{s} = \phi$ (the total magnetic flux inside the circuit).

This corresponds to another experimental situation. By means of a current flowing through a very closely wound cylindrical solenoid of radius R , center at the origin and axis in the z direction, we create a magnetic field, \mathbf{H} , which is essentially confined within the solenoid. However, the vector potential, \mathbf{A} , evidently, cannot be zero everywhere outside the solenoid, because the total flux through every circuit containing the origin is equal to a constant

$$\phi_0 = \int \mathbf{H} \cdot d\mathbf{s} = \int \mathbf{A} \cdot d\mathbf{x}.$$

To demonstrate the effects of the total flux, we begin, as before, with a coherent beam of electrons. (But now there is no need to make wave packets.) The beam is split into two parts, each going on opposite sides of the solenoid, but avoiding it. (The solenoid can be shielded from the electron beam by a thin plate which casts a shadow.) As in the former example, the beams are brought together at F (Fig. 2).

The Hamiltonian for this case is

$$H = \frac{[\mathbf{P} - (e/c)\mathbf{A}]^2}{2m}.$$

In singly connected regions, where $\mathbf{H} = \nabla \times \mathbf{A} = 0$, we can always obtain a solution for the above Hamiltonian by taking $\psi = \psi_0 e^{-iS/\hbar}$, where ψ_0 is the solution when $\mathbf{A} = 0$ and where $\nabla S/\hbar = (e/c)\mathbf{A}$. But, in the experiment discussed above, in which we have a multiply connected region (the region outside the solenoid), $\psi_0 e^{-iS/\hbar}$ is a non-single-valued function¹ and therefore, in general, not a permissible solution of Schrödinger's equation. Nevertheless, in our problem it is still possible to use such solutions because the wave function splits into two parts $\psi = \psi_1 + \psi_2$, where ψ_1 represents the beam on

¹ Unless $\phi_0 = nhc/e$, where n is an integer.

one side of the solenoid¹ and ψ_2 the beam on the opposite side. Each of these beams stays in a simply connected region. We therefore can write

$$\psi_1 = \psi_1^0 e^{-iS_1/\hbar}, \quad \psi_2 = \psi_2^0 e^{-iS_2/\hbar},$$

where S_1 and S_2 are equal to $(e/c) \int A \cdot dx$ along the paths of the first and second beams, respectively. (In Sec. 4, an exact solution for this Hamiltonian will be given, and it will confirm the above results.)

The interference between the two beams will evidently depend on the phase difference,

$$(S_1 - S_2)/\hbar = (e/hc) \int A \cdot dx = (e/hc)\phi_0.$$

This effect will exist, even though there are no magnetic forces acting in the places where the electron beam passes.

In order to avoid fully any possible question of contact of the electron with the magnetic field we note that our result would not be changed if we surrounded the solenoid by a potential barrier that reflects the electrons perfectly. (This, too, is confirmed in Sec. 4.)

It is easy to devise hypothetical experiments in which the vector potential may influence not only the interference pattern but also the momentum. To see this, consider a periodic array of solenoids, each of which is shielded from direct contact with the beam by a small plate. This will be essentially a grating. Consider first the diffraction pattern without the magnetic field, which will have a discrete set of directions of strong constructive interference. The effect of the vector potential will be to produce a shift of the relative phase of the wave function in different elements of the gratings. A corresponding shift will take place in the directions, and therefore the momentum of the diffracted beam.

3. A PRACTICABLE EXPERIMENT TO TEST FOR THE EFFECTS OF A POTENTIAL WHERE THERE ARE NO FIELDS

As yet no direct experiments have been carried out which confirm the effect of potentials where there is no field. It would be interesting therefore to test whether such effects actually exist. Such a test is, in fact, within the range of present possibilities.² Recent experiments^{3,4} have succeeded in obtaining interference from electron beams that have been separated in one case by as much as 0.8 mm.² It is quite possible to wind solenoids which are smaller than this, and therefore to place them between the separate beams. Alternatively, we may obtain localized lines of flux of the right magnitude (the

¹ Dr. Chambers is now making a preliminary experimental study of this question at Bristol.

² L. Marton, Phys. Rev. **85**, 1057 (1952); **90**, 490 (1953). Marton, Simpson, and Suddeth, Rev. Sci. Instr. **25**, 1099 (1954).

³ G. Mollenstedt, Naturwissenschaften **42**, 41 (1955); G. Mollenstedt and H. Düker, Z. Physik **145**, 377 (1956).

magnitude has to be of the order of $\phi_0 = 2\pi ch/e \sim 4 \times 10^{-7}$ gauss cm²) by means of fine permanently magnetized "whiskers".⁵ The solenoid can be used in Marton's device,³ while the whisker is suitable for another experimental setup⁴ where the separation is of the order of microns and the whiskers are even smaller than this.

In principle, we could do the experiment by observing the interference pattern with and without the magnetic flux. But since the main effect of the flux is only to displace the line pattern without changing the interval structure, this would not be a convenient experiment to do. Instead, it would be easier to vary the magnetic flux within the same exposure for the detection of the interference patterns. Such a variation would, according to our previous discussion, alter the sharpness and the general form of the interference bands. This alteration would then constitute a verification of the predicted phenomena.

When the magnetic flux is altered, there will, of course, be an induced electric field outside the solenoid, but the effects of this field can be made negligible. For example, suppose the magnetic flux were suddenly altered in the middle of an exposure. The electric field would then exist only for a very short time, so that only a small part of the beam would be affected by it.

4. EXACT SOLUTION FOR SCATTERING PROBLEMS

We shall now obtain an exact solution for the problem of the scattering of an electron beam by a magnetic field in the limit where the magnetic field region tends to a zero radius, while the total flux remains fixed. This corresponds to the setup described in Sec. 2 and shown in Fig. 2. Only this time we do not split the plane wave into two parts. The wave equation outside the magnetic field region is, in cylindrical coordinates,

$$\left[\frac{\partial^2}{\partial r^2} + \frac{1}{r} \frac{\partial}{\partial r} + \frac{1}{r^2} \left(\frac{\partial}{\partial \theta} + i\alpha \right)^2 + k^2 \right] \psi = 0, \quad (1)$$

where k is the wave vector of the incident particle and $\alpha = -e\phi/ch$. We have again chosen the gauge in which $A_r = 0$ and $A_\theta = \phi/2\pi r$.

The general solution of the above equation is

$$\psi = \sum_{m=-\infty}^{\infty} e^{im\theta} [a_m J_{m+\alpha}(kr) + b_m J_{-(m+\alpha)}(kr)], \quad (2)$$

where a_m and b_m are arbitrary constants and $J_{m+\alpha}(kr)$ is a Bessel function, in general of fractional order (dependent on ϕ). The above solution holds only for $r > R$. For $r < R$ (inside the magnetic field) the solution has been worked out.⁶ By matching the solutions at $r = R$ it is easily shown that only Bessel functions of positive order will remain, when R approaches zero.

⁵ See, for example, Sidney S. Brenner, Acta Met. **4**, 62 (1956).

⁶ L. Page, Phys. Rev. **36**, 444 (1930).

This means that the probability of finding the particle inside the magnetic field region approaches zero with R . It follows that the wave function would not be changed if the electron were kept away from the field by a barrier whose radius also went to zero with R .

The general solution in the limit of R tending to zero is therefore

$$\psi = \sum_{m=-\infty}^{\infty} a_m J_{|m+\alpha|} e^{im\theta}. \quad (3)$$

We must then choose a_m so that ψ represents a beam of electrons that is incident from the right ($\theta=0$). It is important, however, to satisfy the initial condition that the current density,

$$j = \frac{h(\psi^* \nabla \psi - \psi \nabla \psi^*)}{2im} - \frac{e}{mc} A \psi^* \psi, \quad (4)$$

shall be constant and in the x direction. In the gauge that we are using, we easily see that the correct incident wave is $\psi_{\text{inc}} = e^{-ikr} e^{-i\alpha\theta}$. Of course, this wave function holds only to the right of the origin, so that no problem of multiple-valuedness arises.

We shall show in the course of this calculation that the above conditions will be satisfied by choosing $a_m = (-i)^{|m+\alpha|}$, in which case, we shall have

$$\psi = \sum_{m=-\infty}^{\infty} (-i)^{|m+\alpha|} J_{|m+\alpha|} e^{im\theta}.$$

It is convenient to split ψ into the following three parts: $\psi = \psi_1 + \psi_2 + \psi_3$, where

$$\begin{aligned} \psi_1 &= \sum_{m=1}^{\infty} (-i)^{m+\alpha} J_{m+\alpha} e^{im\theta}, \\ \psi_2 &= \sum_{m=-\infty}^{-1} (-i)^{m+\alpha} J_{m+\alpha} e^{im\theta}, \\ &\quad = \sum_{m=1}^{\infty} (-i)^{m-\alpha} J_{m-\alpha} e^{-im\theta}, \quad (5) \\ \psi_3 &= (-i)^{|\alpha|} J_{|\alpha|}. \end{aligned}$$

Now ψ_1 satisfies the simple differential equation

$$\begin{aligned} \frac{\partial \psi_1}{\partial r'} &= \sum_{m=1}^{\infty} (-i)^{m+\alpha} J_{m+\alpha}' e^{im\theta} \\ &= \sum_{m=1}^{\infty} (-i)^{m+\alpha} \frac{J_{m+\alpha-1} - J_{m+\alpha+1}}{2} e^{im\theta}, \quad r' = kr \quad (6) \end{aligned}$$

where we have used the well-known formula for Bessel functions:

$$dJ_{\gamma}(r)/dr = \frac{1}{2}(J_{\gamma-1} - J_{\gamma+1}).$$

As a result, we obtain

$$\begin{aligned} \frac{\partial \psi_1}{\partial r'} &= \frac{1}{2} \sum_{m'=0}^{\infty} (-i)^{m'+\alpha+1} J_{m'+\alpha} e^{i(m'+1)\theta} \\ &\quad - \frac{1}{2} \sum_{m'=2}^{\infty} (-i)^{m'+\alpha-1} J_{m'+\alpha} e^{i(m'-1)\theta} \\ &= \frac{1}{2} \sum_{m'=1}^{\infty} (-i)^{m'+\alpha} J_{m'+\alpha} e^{im'\theta} (-ie^{i\theta} + i^{-1}e^{-i\theta}) \\ &\quad + \frac{1}{2} (-i)^{\alpha} [J_{\alpha+1} - ie^{i\theta} J_{\alpha}]. \end{aligned} \quad (7)$$

So

$$\frac{\partial \psi_1}{\partial r'} = -i \cos \theta \psi_1 + \frac{1}{2} (-i)^{\alpha} (J_{\alpha+1} - ie^{i\theta} J_{\alpha}).$$

This differential equation can be easily integrated to give

$$\psi_1 = A \int_0^{r'} e^{ir' \cos \theta} [J_{\alpha+1} - ie^{i\theta} J_{\alpha}] dr', \quad (8)$$

where

$$A = \frac{1}{2} (-i)^{\alpha} e^{-ir' \cos \theta}.$$

The lower limit of the integration is determined by the requirement that when r' goes to zero, ψ_1 also goes to zero because, as we have seen, ψ_1 includes Bessel functions of positive order only.

In order to discuss the asymptotic behavior of ψ_1 , let us write it as $\psi_1 = A [I_1 - I_2]$, where

$$\begin{aligned} I_1 &= \int_0^{\infty} e^{ir' \cos \theta} [J_{\alpha+1} - ie^{i\theta} J_{\alpha}] dr', \\ I_2 &= \int_r^{\infty} e^{ir' \cos \theta} [J_{\alpha+1} - ie^{i\theta} J_{\alpha}] dr'. \end{aligned} \quad (9)$$

The first of these integrals is known⁷:

$$\int_0^{\infty} e^{i\beta r} J_{\alpha}(kr) = \frac{e^{i[\alpha \arcsin(\beta/k)]}}{(k^2 - \beta^2)^{\frac{1}{2}}}, \quad 0 < \beta < k, \quad -2 < \alpha$$

In our cases, $\beta = \cos \theta$, $k = 1$, so that

$$I_1 = \left[\frac{e^{i\alpha(\frac{1}{2}\pi - |\theta|)}}{|\sin \theta|} - ie^{i\theta} \frac{e^{i(\alpha+1)(\frac{1}{2}\pi - |\theta|)}}{|\sin \theta|} \right]. \quad (10)$$

Because the integrand is even in θ , we have written the final expression for the above integral as a function of $|\theta|$ and of $|\sin \theta|$. Hence

$$\begin{aligned} I_1 &= e^{i\alpha(\frac{1}{2}\pi - |\theta|)} \left[\frac{ie^{-i|\theta|} - ie^{i\theta}}{|\sin \theta|} \right] \\ &= 0 \quad \text{for } \theta < 0, \\ &= e^{-i\alpha\theta} 2i^{\alpha} \quad \text{for } \theta > 0, \end{aligned} \quad (11)$$

where we have taken θ as going from $-\pi$ to π .

⁷ See, for example, W. Gröbner and N. Hofreiter, *Integraltafel* (Springer-Verlag, Berlin, 1949).

We shall see presently that I_1 represents the largest term in the asymptotic expansion of ψ_1 . The fact that it is zero for $\theta < 0$ shows that this part of ψ_1 passes (asymptotically) only on the upper side of the singularity. To explain this, we note that ψ_1 contains only positive values of m , and therefore of the angular momentum. It is quite natural then that this part of ψ_1 goes on the upper side of the singularity. Similarly, since according to (5)

$$\psi_2(r', \theta, \alpha) = \psi_1(r', -\theta, -\alpha),$$

it follows that ψ_2 will behave oppositely to ψ_1 in this regard, so that together they will make up the correct incident wave.

Now, in the limit of $r' \rightarrow \infty$ we are allowed to take in the integrand of I_2 the first asymptotic term of J_α ,⁸ namely $J_\alpha \rightarrow (2/\pi r')^{\frac{1}{2}} \cos(r' - \frac{1}{2}\alpha - \frac{1}{4}\pi)$. We obtain

$$I_2 = \int_r^\infty e^{ir' \cos\theta} (J_{\alpha+1} - ie^{i\theta} J_\alpha) dr' \rightarrow C + D, \quad (12)$$

where

$$C = \int_r^\infty e^{ir' \cos\theta} [\cos(r' - \frac{1}{2}(\alpha+1)\pi - \frac{1}{4}\pi)] \frac{dr'}{(r')^{\frac{1}{2}}} \left(\frac{2}{\pi} \right)^{\frac{1}{2}}, \quad (13)$$

$$D = \int_r^\infty e^{ir' \cos\theta} [\cos(r' - \frac{1}{2}\alpha - \frac{1}{4}\pi)] \frac{dr'}{(r')^{\frac{1}{2}}} \left(\frac{2}{\pi} \right)^{\frac{1}{2}} (-i)e^{i\theta}.$$

Then

$$\begin{aligned} C &= \int_r^\infty e^{ir' \cos\theta} [e^{i(r' - \frac{1}{2}(\alpha+1)\pi - \frac{1}{4}\pi)} \\ &\quad + e^{-i(r' - \frac{1}{2}(\alpha+1)\pi - \frac{1}{4}\pi)}] \frac{dr'}{(2\pi r')^{\frac{1}{2}}} \\ &= \left(\frac{2}{\pi} \right)^{\frac{1}{2}} \frac{(-i)^{\alpha+1}}{(1+\cos\theta)^{\frac{1}{2}}} \int_{[r'(1+\cos\theta)]^{\frac{1}{2}}}^\infty \exp(+iz^2) dz \\ &\quad + \left(\frac{2}{\pi} \right)^{\frac{1}{2}} \frac{i^{\alpha+1}}{(1-\cos\theta)^{\frac{1}{2}}} \int_{[r'(1-\cos\theta)]^{\frac{1}{2}}}^\infty \exp(-iz^2) dz, \end{aligned} \quad (14)$$

where we have put

$$z = [r'(1+\cos\theta)]^{\frac{1}{2}} \quad \text{and} \quad z = [r'(1-\cos\theta)]^{\frac{1}{2}},$$

respectively.

Using now the well-known asymptotic behavior of the error function,⁹

$$\begin{aligned} \int_a^\infty \exp(iz^2) dz &\rightarrow -\frac{i \exp(ia^2)}{2}, \\ \int_a^\infty \exp(-iz^2) dz &\rightarrow \frac{-i \exp(-ia^2)}{2}, \end{aligned} \quad (15)$$

⁸ E. Jahnke and F. Emde, *Tables of Functions* (Dover Publications, Inc., New York, 1943), fourth edition, p. 138.

⁹ Reference 8, p. 24.

we finally obtain

$$C = \left[\frac{(-i)^{\alpha+1}}{(2\pi)^{\frac{1}{2}}} \frac{e^{ir'}}{[r'(1+\cos\theta)^2]^{\frac{1}{2}}} \right. \\ \left. + \frac{i^{\alpha+1}}{(2\pi)^{\frac{1}{2}}} \frac{e^{-ir'}}{[r'(1-\cos\theta)^2]^{\frac{1}{2}}} \right] e^{ir' \cos\theta}, \quad (16)$$

$$D = \left[\frac{(-i)^{\alpha+1}}{(2\pi)^{\frac{1}{2}}} \frac{e^{ir'}}{[r'(1+\cos\theta)^2]^{\frac{1}{2}}} \right. \\ \left. + \frac{i^{\alpha+1}}{(2\pi)^{\frac{1}{2}}} \frac{e^{-ir'}}{[r'(1-\cos\theta)^2]^{\frac{1}{2}}} \right] e^{ir' \cos\theta} (-i)e^{i\theta}. \quad (17)$$

Now adding (16) and (17) together and using (13) and (9), we find that the term of $1/(r')^{\frac{1}{2}}$ in the asymptotic expansion of ψ_1 is

$$\frac{(-i)^{\frac{1}{2}}}{2(2\pi)^{\frac{1}{2}}} \left[(-1)^\alpha \frac{e^{ir'}}{(r')^{\frac{1}{2}}} \frac{1+e^{i\theta}}{1+\cos\theta} + i \frac{e^{-ir'}}{(r')^{\frac{1}{2}}} \frac{1-e^{i\theta}}{1-\cos\theta} \right]. \quad (18)$$

Using again the relation between ψ_1 and ψ_2 we obtain for the corresponding term in ψ_2

$$\frac{(-i)^{\frac{1}{2}}}{2(2\pi)^{\frac{1}{2}}} \left[(-1)^{-\alpha} \frac{e^{ir'}}{(r')^{\frac{1}{2}}} \frac{1+e^{-i\theta}}{1+\cos\theta} + i \frac{e^{-ir'}}{(r')^{\frac{1}{2}}} \frac{1-e^{-i\theta}}{1-\cos\theta} \right]. \quad (19)$$

Adding (18) and (19) and using (11), we finally get

$$\begin{aligned} \psi_1 + \psi_2 &\rightarrow \frac{(-i)^{\frac{1}{2}}}{(2\pi)^{\frac{1}{2}}} \left[\frac{ie^{-ir'}}{(r')^{\frac{1}{2}}} + \frac{e^{ir'}}{(r')^{\frac{1}{2}}} \frac{\cos(\pi\alpha - \frac{1}{2}\theta)}{\cos(\frac{1}{2}\theta)} \right] \\ &\quad + e^{-i(r' \cos\theta + \alpha\theta)}. \end{aligned} \quad (20)$$

There remains the contribution of ψ_3 , whose asymptotic behavior is [see Eq. (12)]

$$(-i)^{|\alpha|} J_{|\alpha|}(r') \rightarrow (-i)^{|\alpha|} \left(\frac{2}{\pi r'} \right)^{\frac{1}{2}} \cos(r' - \frac{1}{4}\pi - \frac{1}{2}|\alpha|\pi).$$

Collecting all terms, we find

$$\psi = \psi_1 + \psi_2 + \psi_3 \rightarrow e^{-i(\alpha\theta + r' \cos\theta)} + \frac{e^{ir'}}{(2\pi ir')^{\frac{1}{2}}} \frac{\sin(\pi\alpha - \theta/2)}{\cos(\theta/2)}, \quad (21)$$

where the \pm sign is chosen according to the sign of α .

The first term in equation (21) represents the incident wave, and the second the scattered wave.¹⁰ The scattering cross section is therefore

$$\sigma = \frac{\sin^2 \pi\alpha}{2\pi} \frac{1}{\cos^2(\theta/2)}. \quad (22)$$

¹⁰ In this way, we verify, of course, that our choice of the a_m for Eq. (3) satisfies the correct boundary conditions.

When $\alpha = n$, where n is an integer, then σ vanishes. This is analogous to the Ramsauer effect.¹¹ σ has a maximum when $\alpha = n + \frac{1}{2}$.

The asymptotic formula (21) holds only when we are not on the line $\theta = \pi$. The exact solution, which is needed on this line, would show that the second term will combine with the first to make a single-valued wave function, despite the non-single-valued character of the two parts, in the neighborhood of $\theta = \pi$. We shall see this in more detail presently for the special case $\alpha = n + \frac{1}{2}$.

In the interference experiment discussed in Sec. 2, diffraction effects, represented in Eq. (21) by the scattered wave, have been neglected. Therefore, in this problem, it is adequate to use the first term of Eq. (21). Here, we see that the phase of the wave function has a different value depending on whether we approach the line $\theta = \pm\pi$ from positive or negative angles, i.e., from the upper or lower side. This confirms the conclusions obtained in the approximate treatment of Sec. 2.

We shall discuss now the two special cases that can be solved exactly. The first is the case where $\alpha = n$. Here, the wave function is $\psi = e^{-ikr} e^{-ia\theta}$, which is evidently single-valued when α is an integer. (It can be seen by direct differentiation that this is a solution.)

The second case is that of $\alpha = n + \frac{1}{2}$. Because $J_{(n+\frac{1}{2})}(r)$ is a closed trigonometric function, the integrals for ψ can be carried out exactly.

The result is

$$\psi = \frac{i^{\frac{1}{2}}}{\sqrt{2}} e^{-i(\theta + r' \cos \theta)} \int_0^{[r'(1+\cos \theta)]^{\frac{1}{2}}} \exp(iz^2) dz. \quad (23)$$

This function vanishes on the line $\theta = \pi$. It can be seen that its asymptotic behavior is the same as that of Eq. (2) with α set equal to $n + \frac{1}{2}$. In this case, the single-valuedness of ψ is evident. In general, however, the behavior of ψ is not so simple, since ψ does not become zero on the line $\theta = \pi$.

5. DISCUSSION OF SIGNIFICANCE OF RESULTS

The essential result of the previous discussion is that in quantum theory, an electron (for example) can be influenced by the potentials even if all the field regions are excluded from it. In other words, in a field-free multiply-connected region of space, the physical properties of the system still depend on the potentials.

It is true that all these effects of the potentials depend only on the gauge-invariant quantity $\oint \mathbf{A} \cdot d\mathbf{x} = \oint \mathbf{H} \cdot d\mathbf{s}$, so that in reality they can be expressed in terms of the fields inside the circuit. However, according to current relativistic notions, all fields must interact only locally. And since the electrons cannot reach the regions where the fields are, we cannot interpret such effects as due to the fields themselves.

¹¹ See, for example, D. Bohm, *Quantum Theory* (Prentice-Hall, Inc., Englewood Cliffs, New Jersey, 1951).

In classical mechanics, we recall that potentials cannot have such significance because the equation of motion involves only the field quantities themselves. For this reason, the potentials have been regarded as purely mathematical auxiliaries, while only the field quantities were thought to have a direct physical meaning.

In quantum mechanics, the essential difference is that the equations of motion of a particle are replaced by the Schrödinger equation for a wave. This Schrödinger equation is obtained from a canonical formalism, which cannot be expressed in terms of the fields alone, but which also requires the potentials. Indeed, the potentials play a role in Schrödinger's equation, which is analogous to that of the index of refraction in optics. The Lorentz force [$e\mathbf{E} + (e/c)\mathbf{v} \times \mathbf{H}$] does not appear anywhere in the fundamental theory, but appears only as an approximation holding in the classical limit. It would therefore seem natural at this point to propose that, in quantum mechanics, the fundamental physical entities are the potentials, while the fields are derived from them by differentiations.

The main objection that could be raised against the above suggestion is grounded in the gauge invariance of the theory. In other words, if the potentials are subject to the transformation $A_\mu \rightarrow A'_\mu = A_\mu + \partial\psi/\partial x_\mu$, where ψ is a continuous scalar function, then all the known physical quantities are left unchanged. As a result, the same physical behavior is obtained from any two potentials, $A_\mu(x)$ and $A'_\mu(x)$, related by the above transformation. This means that insofar as the potentials are richer in properties than the fields, there is no way to reveal this additional richness. It was therefore concluded that the potentials cannot have any meaning, except insofar as they are used mathematically, to calculate the fields.

We have seen from the examples described in this paper that the above point of view cannot be maintained for the general case. Of course, our discussion does not bring into question the gauge invariance of the theory. But it does show that in a theory involving only local interactions (e.g., Schrödinger's or Dirac's equation, and current quantum-mechanical field theories), the potentials must, in certain cases, be considered as physically effective, even when there are no fields acting on the charged particles.

The above discussion suggests that some further development of the theory is needed. Two possible directions are clear. First, we may try to formulate a nonlocal theory in which, for example, the electron could interact with a field that was a finite distance away. Then there would be no trouble in interpreting these results, but, as is well known, there are severe difficulties in the way of doing this. Secondly, we may retain the present local theory and, instead, we may try to give a further new interpretation to the poten-

tials. In other words, we are led to regard $A_\mu(x)$ as a physical variable. This means that we must be able to define the physical difference between two quantum states which differ only by gauge transformation. It will be shown in a future paper that in a system containing an undefined number of charged particles (i.e., a superposition of states of different total charge), a new Hermitian operator, essentially an angle variable, can be introduced, which is conjugate to the charge density and which may give a meaning to the gauge. Such states have actually been used in connection with

recent theories of superconductivity and superfluidity¹² and we shall show their relation to this problem in more detail.

ACKNOWLEDGMENTS

We are indebted to Professor M. H. L. Pryce for many helpful discussions. We wish to thank Dr. Chambers for many discussions connected with the experimental side of the problem.

¹² See, for example, C. G. Kuper, *Advances in Physics*, edited by N. F. Mott (Taylor and Francis, Ltd., London, 1959), Vol. 8, p. 25, Sec. 3, Par. 3.

Chapter 3

FOUNDATIONS

- [3.1] M. V. Berry, "Quantal Phase Factors Accompanying Adiabatic Changes," *Proc. R. Lond. A* **392** (1984) 45–57 124
- [3.2] B. Simon, "Holonomy, the Quantum Adiabatic Theorem, and Berry's Phase," *Phys. Rev. Lett.* **51** (1983) 2167–2170 137
- [3.3] F. Wilczek and A. Zee, "Appearance of Gauge Structure in Simple Dynamical Systems," *Phys. Rev. Lett.* **52** (1984) 2111–2114 141
- [3.4] Y. Aharonov and J. Anandan, "Phase Change during a Cyclic Quantum Evolution," *Phys. Rev. Lett.* **58** (1987) 1593–1596 145
- [3.5] J. Samuel and R. Bhandari, "General Setting for Berry's Phase," *Phys. Rev. Lett.* **60** (1988) 2339–2342 149
- [3.6] H. Kuratsuji and S. Iida, "Effective Action for Adiabatic Process," *Prog. Theo. Phys.* **74** (1985) 439–445 153
- [3.7] J. Moody, A. Shapere and F. Wilczek, "Adiabatic Effective Lagrangians" * 160

* Original Contribution.

3

Foundations

We saw in the previous chapter how molecular and optical physicists originally learned to take account of geometric phases. In optics, Pancharatnam's phase led to measurable interference effects. In molecular physics, the molecular Aharonov–Bohm effect was found to have a significant effect on molecular dynamics near to a degeneracy of two electronic energy levels. It was left for Berry to fully appreciate the universal significance of these phases, to show that whenever an adiabatic approximation applies, we may expect to find a geometric phase. The crucial role of adiabaticity is to make sure that the cyclic variation of parameters leads to cyclic evolution. (If one is otherwise able to ensure a cyclic evolution, then as Aharonov and Anandan have pointed out, the adiabatic hypothesis may be dispensed with.) It is the ubiquity of the adiabatic approximation, in both theoretical and experimental physics, that has led to the wide application of Berry's basic observation.

We now come to Berry's original paper on the quantal adiabatic phase [3.1], and its subsequent elaborations and generalizations. Among the latter are Wilczek and Zee's work on non-Abelian phases associated with adiabatic evolution of degenerate hamiltonians [3.3], and Aharonov and Anandan's geometric phase for cyclic evolution, not necessarily adiabatic, in projective Hilbert space [3.4]. The final entry of this chapter is a detailed discussion of adiabatic phases and related issues in the context of the Born–Oppenheimer approximation [3.7].

The paper “Quantal phase factors accompanying adiabatic changes” [3.1], elegantly presents the key concepts surrounding what we have come to know as Berry's phase—the gauge invariance of the adiabatic phase, the expression of the phase as the integral of a two-form over an enclosed surface, the theorem on degeneracies of a complex hamiltonian. The example of a spin in a slowly changing magnetic field, which has become a paradigm in studies (and measurements) of geometric phases, is shown to lead to magnetic-monopole-like effects. Berry proposes an experiment to measure the phase, by splitting a beam of coherently polarized electrons in a magnetic field, rotating the field around one of the beams, and measuring the

resulting phase difference by interference.

In the next paper [3.2], Simon gives a mathematical interpretation of Berry's phase as the holonomy of a complex line bundle. To make his paper more accessible to non-mathematicians, we include an appendix on holonomy and fiber bundles at the end of this introduction. The appendix may also help to clarify mathematical aspects of several other papers in this and later chapters.

Berry's paper assumes that the energy eigenstate undergoing adiabatic evolution is non-degenerate. If this is so, then a cyclic variation C of the external parameters will return the system to its original state, multiplied by a complex number of unit modulus, the product of a dynamical phase and a geometric phase $\exp i\gamma(C)$:

$$|\Psi(T)\rangle = \exp \left\{ -\frac{i}{\hbar} \int^T E(t) dt \right\} \cdot \exp i\gamma(C) |\Psi(0)\rangle \quad (3.1)$$

However, when the state is degenerate over the full course of its evolution, the system need not return to the original eigenstate, but only to one of the degenerate states. As observed in the paper by Wilczek and Zee [3.3], the accumulated "phase" for an N -fold degenerate level will actually be a $U(N)$ matrix in general:

$$|\Psi_\alpha(T)\rangle = \exp \left\{ -\frac{i}{\hbar} \int^T E(t) dt \right\} U_{\alpha\beta} |\Psi_\beta(0)\rangle \quad (3.2)$$

where α runs from 1 to N . The unitary matrix $U_{\alpha\beta}$ may be expressed as a time-ordered exponential integral:

$$U = T \exp \int_0^T A(t) dt; \quad A_{\alpha\beta} \equiv \langle \Psi_\alpha | \dot{\Psi}_\beta \rangle \quad (3.3)$$

Wilczek and Zee present several examples of systems exhibiting these non-Abelian phase factors; many more may be found in Chapter 4. Systems with global degeneracies are quite interesting experimentally, because the $U(N)$ adiabatic phase represents the leading contribution to mixings between degenerate levels.

Another important generalization of Berry's phase, one which has recently provoked a great deal of interest, is the phase of Aharonov and Anandan [3.4]. Recall that Berry studied cyclic evolution in parameter space, which in the adiabatic limit always leads to cyclic evolution in the projective Hilbert space \mathcal{R} of states $|\Psi\rangle$ modulo phases, as shown by Eq. (3.1). Thus the essential ingredient needed in order to define a geometric phase is the closed path in \mathcal{R} ; adiabaticity is the sufficient (but not necessary) condition which guarantees the existence of such a path. Aharonov and Anandan take \mathcal{R} as their starting point, and show how any closed path in

\mathcal{R} has a geometric phase associated to it in a natural way. The underlying parameter space plays no fundamental role in their considerations, although of course the Aharonov–Anandan (AA) phase reduces to Berry’s phase when the closed path arises from adiabatic evolution of external parameters. In their paper, they discuss three examples of cyclic evolution which are not necessarily adiabatic. One example, an electron in an external electromagnetic field, is used to demonstrate that Aharonov–Bohm effect is a special case of the AA phase. The other applications concern the much-studied system of a spin precessing in a magnetic field—in fact, the AA phase has been directly measured in magnetic resonance experiments [4.6]. In general, however, it seems to be rather difficult to ensure cyclic evolution in projective Hilbert space in the absence of an adiabatic hypothesis.

The paper [3.5] by Samuel and Bhandari shows that even the hypothesis of cyclic evolution may be dispensed with. Therein, these authors propose a further generalization of both the Berry and AA phases. Harking back to work of Pancharatnam in optics [2.2] they point out that there is a natural way to close an open path, if one is given a metric on the underlying space. Namely, one simply joins the two ends by a geodesic. So if the state of a system undergoes a non-cyclic evolution as its parameters are varied over a time interval from 0 to T , we can “close” its evolution with a geodesic (in projective Hilbert space), and obtain the corresponding geometric phase. This shows that, provided we are given a metric on the projective Hilbert space \mathcal{R} , and *modulo* questions of uniqueness, there is a geometric phase associated to every path, closed or not! Now over Hilbert space there is indeed a natural *gauge invariant* metric, which Berry has discussed in the introduction to this volume [1.1]. If $|dn\rangle$ is a tangent vector to \mathcal{R} at $|n\rangle$, then the metric tensor is

$$ds^2 = \langle \partial_i n | (1 - |n\rangle\langle n|) |\partial_j n \rangle dX_i dX_j = g_{ij} dX_i dX_j \quad (3.4)$$

This metric is not positive definite, and does not in general determine a unique geodesic between any two points. However, as Samuel and Bhandari show, the geometric phase obtained is *independent* of the choice of geodesic. It has a nice physical interpretation, too: when the initial and final states are “in phase” (so that the geometric phase is 1), then the norm of the sum of the two states reaches a maximum. The only condition for the relative phase to be well-defined is that the initial and final states have a nonzero overlap. All this is strongly reminiscent of Pancharatnam’s connection, which allows one to compare relative phases of beams of light in very different polarization states, not just nearby ones. There, the relative phase of two states is defined by parallel transporting one of the states along a geodesic on the Poincaré sphere connecting the two. Just as in the case at hand, Pancharatnam’s condition for two beams of light to be in phase implies that the intensity of their sum is maximal. It is hard to imagine anything more general than the geometric phase of Samuel and Bhandari, which applies to essentially any type of quantum evolution imaginable.

The remaining two entries in this chapter treat adiabatic phases in the context of the Born-Oppenheimer approximation. Kuratsuji and Iida were, as far as we know, the first to discuss Berry's phase in a path-integral framework, as opposed to a Schrödinger description [3.6]. Such a framework lends itself well to the study of Born-Oppenheimer systems, where one is interested in separating the integrations over "fast" electronic variables and "slow" nuclear variables. Typically, one performs the electronic path integration first, in a fixed nuclear background, and making the adiabatic approximation that electronic transitions do not occur. The result is an effective action involving only the nuclear coordinates. (This procedure is discussed in detail in the final paper, by Moody, Shapere, and Wilczek [3.7].) Now when, for example, the nuclei traverse a closed path, at a typically slow nuclear velocity, the electronic wavefunction will acquire an adiabatic phase. Berry's phase must be taken into account in the nuclear path integral. As Kuratsuji and Iida point out, the appropriate modification consists of a simple addition to the nuclear effective action.

The path-integral description also makes it easy to address questions about the semiclassical limit of nuclear motion. Here Berry's phase has a strikingly direct effect on the energy levels. If C is a closed classical path and $(P(t), Q(t))$ the corresponding trajectory in phase space, then Kuratsuji and Iida, and also Wilkinson [6.6], find the following quantization rule:

$$\oint P \cdot dX = \left(n + \frac{\alpha}{4} - \frac{\Gamma(C)}{2\pi} \right) 2\pi\hbar. \quad (3.5)$$

Here $\Gamma(C)$ is the Berry's phase associated with the path and α is an integer known as the Keller-Maslov index, and determined by continuing the WKB wavefunction around the turning points.

"Adiabatic Effective Lagrangians" [3.7] is the paper referred to in [4.3] as Ref. 7, which has not appeared in print until now. Therein, the general procedure of "integrating out" fast degrees of freedom is explained, emphasizing its application to molecular systems. The Born-Oppenheimer method in both its Hamiltonian and Lagrangian forms is laid out in detail, taking into account the possibility of an adiabatic phase. As we have said, the phase enters into the nuclear effective theories as a canonically coupled background gauge potential. Perhaps surprisingly, the Born-Oppenheimer Hamiltonian and Lagrangian are not related by a Legendre transformation. An explicit expansion giving corrections to adiabatic evolution is derived, and used to give a proof of the adiabatic theorem. For smooth classical motions of the nuclei, tunneling corrections are exponentially suppressed, away from level crossings. On the other hand, corrections to phase evolution are only suppressed by powers of the expansion parameter. These may be accounted for directly by adding higher-derivative terms to the nuclear effective Lagrangian.

Appendix: Monopoles, Holonomy and Fiber Bundles

A “geometric phase” is what mathematicians would call a $U(1)$ holonomy, and the natural mathematical context for holonomy is the theory of fiber bundles. Because they will continue to arise implicitly and explicitly throughout this book, we would now like to give a brief introduction to fiber bundles and some related mathematical concepts, such as connections and parallel transport. Wherever possible, we shall try to translate mathematics into physics, and to illustrate key concepts by way of a particular physical example, the magnetic monopole.

Before saying what a fiber bundle is, let us jump ahead to explain what a fiber bundle is good for. Suppose the state of a system is described by two sets of parameters, which we shall refer to as internal and external. Suppose further that the internal parameters change in a well-defined way when we vary the external parameters. Thus, in quantum mechanics with a Hamiltonian depending on slowly-varying external parameters $R(t)$, we may think of the phase of the wavefunction as an internal parameter that depends on $R(t)$ through the energies and through the adiabatic phase

$$\exp \int A_n(R) \cdot dR \quad (3.A.1)$$

where $A_n(R) \equiv \langle n(R) | \nabla | n(R) \rangle$. (To avoid irrelevant complications, we suppose the energy $E_n(R) = 0$.) Where does the evolving wavefunction live? It is not enough to say it is a function of R , since its phase is a function of the whole previous history. Rather, we may think of it as living on a space that looks at least locally like $\{R\} \times U(1)$, and whose global topological properties are encoded in $A_n(R)$. In fact, globally this product may be “twisted”—this will be the case if $A_n(R)$ has unremovable singularities, like the Dirac string of a magnetic monopole gauge potential. To describe such twisted products, in a way that avoids talking about singularities, we introduce the following construction.

A bundle E with fiber F is a generalization of the direct product of two spaces $M \times F$. Locally, in a small neighborhood U of any point of M , the bundle looks like $U \times F$, but globally, it may be topologically twisted. M is called the base space, and the fibers of E (which are all homeomorphic to F) may conveniently be thought of as residing vertically over the points of M , with one fiber above each point (see Fig.A.1).

The Möbius strip is a simple example of a fiber bundle. As depicted in Fig. A.2, its base space is a circle and its fiber is a line segment, which we may take to be the interval $[-1, 1]$. Every point on the circle has a neighborhood U over which the bundle is homeomorphic to $U \times [-1, 1]$, but globally there is a twist. One way to see the twist is to choose a homeomorphism i between the fiber over x_0 and $[-1, 1]$, and to try to extend it continuously around the circle. When we come back to x_0 , we find that its “direction” has been

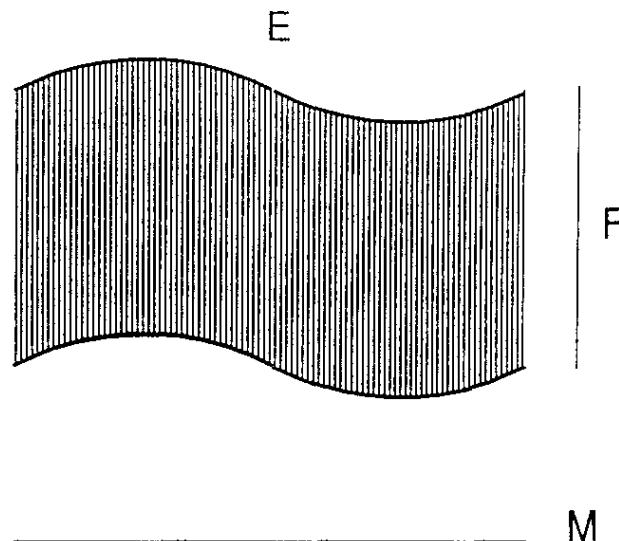


Figure A.1 A fiber bundle with total space E , fiber F , and base space M .

reversed: topologically, transporting i around the circle gives back $-i$. This minus sign is the simplest example of a non-trivial holonomy. It is essentially the same minus sign found by Herzberg and Longuet-Higgins in [2.3], that arose from transporting an eigenstate of a real Hamiltonian around a conical degeneracy.

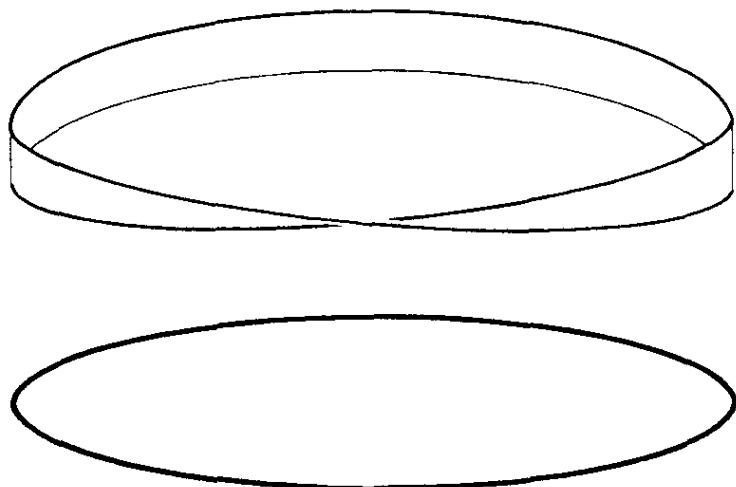


Figure A.2 The Möbius strip is one of the simplest nontrivial examples of a fiber bundle.

We shall be particularly interested in bundles that have a symmetry group G acting upon the fiber. These are properly known as G -bundles, or bundles with structure group G . When the fiber F is the group G itself, we say that we have a principal G -bundle. The action of G on the fibers is to be thought of as a change of coordinates for each fiber. For example, the fiber may be an n -dimensional complex vector space. To label a vector, we must choose a basis; but there remains an action of $U(n)$ representing our freedom to change bases. Suppose now that we cover M with a collection

of small neighborhoods $\{U_i\}$, that have intersections $U_i \cap U_j \equiv U_{ij}$, and over each of which the bundle is isomorphic to $U_i \times F$. Then to sew all the mini-bundles $U_i \times F$ together to make the full bundle E , we need rules for how to identify elements of overlapping mini-bundles (*i.e.*, how to relate two different coordinate systems on the fibers residing over U_{ij}). These rules are known as *transition functions*, and are responsible for all the global topological twisting. They are maps g_{ij} from U_{ij} into the structure group G , and may be regarded as taking elements (p, x) of $U_i \times F$ to (p, gx) in $U_j \times F$, for $p \in U_{ij}$.

Let us illustrate the above abstractions with the “original” example of a fiber bundle, known to physicists as the magnetic monopole¹ and to mathematicians as the Hopf bundle.² This is a bundle over the two-dimensional sphere, with fiber $U(1)$. Physically, it is related to the interaction of an electric charge moving on the surface of the sphere with a magnetic monopole charge at the center of the sphere. Now if the electric charge moves very slowly, then the classical force between the two charges will be vanishingly small. But, of course, classical electromagnetism is not the whole story. Quantum mechanically, there is a very significant interaction, found by Aharonov and Bohm [2.6], involving a geometric phase that depends only on the shape of the path taken by the charge, and not on how fast that path is traversed. The interaction may be described as follows. The magnetic charge has a field $\mathbf{B} = g \hat{\mathbf{r}}/r^2$, and transporting the charge around a closed loop C produces a phase change of the charged particle’s wavefunction of $e^{i\Phi(C)}$, where $\Phi(C)$ is the enclosed magnetic flux (if the electric charge is equal to 1). In the absence of magnetic monopoles, $\nabla \cdot \mathbf{B} = 0$, so we can write \mathbf{B} as the curl of a gauge potential \mathbf{A} (also known as a vector potential or a connection) and express the enclosed flux as

$$\Phi = \int_S \mathbf{B} \cdot d\mathbf{S} = \oint_C \mathbf{A} \cdot d\mathbf{l} \quad (3.A.2)$$

The right-hand integral shows clearly that the accumulated phase depends only on the geometry of the path C . (Another advantage of using a gauge potential is that it is usually easier to do a line integral than a surface integral.)

When magnetic monopoles are present,

$$\nabla \cdot \mathbf{B} = \sum_{\text{monopoles}} g_i \delta^{(3)}(\mathbf{x} - \mathbf{x}_i) \quad (3.A.3)$$

where g_i is the magnetic charge at \mathbf{x}_i , and the gauge potential A will in general be singular. For example, for a single monopole at the origin, a common choice of gauge potential is

$$A_\phi = g(1 - \cos \theta) \quad A_r = A_\theta = 0 \quad (3.A.4)$$

in spherical coordinates. Note that A_ϕ has a “Dirac string” singularity along the line $\theta = \pi$: there the magnitude of A_ϕ ,

$$g^{\phi\phi} A_\phi A_\phi = \frac{g^2(1 - \cos\theta)^2}{\sin^2\theta},$$

blows up. The string can be moved around by means of a gauge transformation $\mathbf{A} \rightarrow \mathbf{A} + \nabla\Lambda$, but it cannot be removed.

In order to avoid singular gauge potentials, Wu and Yang³ introduced the following construction. They covered S^2 with two patches S^+ and S^- ,

$$\begin{aligned} S^+ &\equiv \{(\theta, \phi) : 0 \leq \theta < \frac{\pi}{2} + \epsilon\} \\ S^- &\equiv \{(\theta, \phi) : \frac{\pi}{2} - \epsilon < \theta \leq \pi\}, \end{aligned} \quad (3.A.5)$$

two open sets that respectively contain the northern and southern hemispheres, and whose intersection is an open set containing the equator. Over each patch, the restricted bundle is isomorphic to the trivial bundle $S^\pm \times U(1)$, with elements $(\theta, \phi, e^{i\alpha_\pm})$. The transition function g_{+-} that connects the upper and lower halves will in general be of the form

$$g_{+-} : (\theta, \phi, e^{i\alpha_-}) \longrightarrow (\theta, \phi, e^{i f_{+-}(\theta, \phi)} \cdot e^{i\alpha_+}) \quad (3.A.6)$$

where θ and ϕ are in $S^+ \cap S^-$. Obviously, the winding number of f_{+-} must be an integer. This integer does not depend on our patching scheme, and in fact it is a topological invariant that gives a complete topological characterization of all principal $U(1)$ bundles over the two-sphere. We shall see that this integer is proportional to the magnetic charge of the monopole.

Using the above construction, we can specify a non-singular gauge potential for the magnetic monopole field. One choice is

$$\begin{aligned} A_\phi^+ &= g(\cos\theta + 1) && \text{over } S^+ \\ A_\phi^- &= g(\cos\theta - 1) && \text{over } S^- \end{aligned} \quad (3.A.7)$$

These are each nonsingular in their respective domains, give the right curl, and are related on the equator $S^+ \cap S^-$ by a gauge transformation

$$A_\phi^+ = A_\phi^- + \partial_\phi(2g\phi). \quad (3.A.8)$$

This gauge transformation corresponds to having for a transition function $f_{+-}(\theta, \phi) = 2g\phi$. One may show that the flux Φ of Eq.(3.A.2), a physical quantity, is invariant under this transformation, and that it is equal to the monopole charge times the solid angle subtended by the closed path C . So as claimed, $2g$ is equal to the winding number of f_{+-} , and g must be a multiple of $1/2$. This is Dirac’s famous quantization condition.

Let us now look at the charge-monopole interaction in a more general setting. To repeat, the effect of this interaction on a unit charge that is slowly transported through the field of a magnetic monopole along any path C , open or closed, is to multiply the wavefunction by a geometric phase

$$\exp i \int_C A \cdot dl \quad (3.A.9)$$

The gauge potential thus gives us a rule for lifting a curve in the charged particle's position space to a curve in the $U(1)$ bundle of position wavefunctions. This is an example of the notion of parallel transport, to which we now turn.

In general, a rule for parallel transport is a rule for changing ones fiber coordinate as one moves horizontally over the base space M . It describes a way to lift a curve in the base manifold to a curve in the bundle. This is often necessary for performing explicit computations, and it gives us valuable geometric information about the bundle, as well. It is usually most convenient to define such a rule infinitesimally, that is, as a rule for lifting vectors from M to E . If the vectors lifted are the tangent vectors of a curve C in M , their liftings may then be integrated to give a lifting of C . So let v be a vector tangent to M , and let A be a linear map taking v to a vector tangent to the fiber G , which we denote by A_v . A now specifies a rule for lifting vectors,

$$v \longrightarrow v + A_v \quad (3.A.10)$$

Locally, over an open set U where the bundle looks like $U \times F$, (v, A_v) is a vector in the tangent space of $U \times F \subseteq E$. Thus, A_v is the "vertical" component of the lifted vector. A is known to mathematicians as a connection, and is mathematically identical to a gauge potential.

There is a certain amount of freedom in choosing a connection, corresponding to our freedom to choose coordinates for the fibers. A different choice of a homeomorphism between the fiber above x and the fibering space F yields a different connection. Given one rule for lifting a curve, we obtain another lifting from M to E by multiplying each point of the lifted curve by a different group element $g(x)$. This amounts to rotating the coordinates of the fiber by $g(x)$. Geometrically, the new lifting is equivalent to the old one; we have just chosen to describe it using different coordinates. Although the choice we make is rather arbitrary, it is necessary to make one, in order to specify a lifting. Physicists call this a choice of gauge, and different choices are related by gauge transformations. Under such a transformation, A transforms just as a gauge potential should:

$$A_v(x) \longrightarrow g(x)A_v(x)g(x)^{-1} + g(x)\nabla_v g(x)^{-1} \quad (3.A.11)$$

We found an example of this type of transformation in Eq.(3.A.8), relating A^+ and A^- at the equator.

Parallel transport does not necessarily lift a closed curve C in M to a closed curve in the total space of the bundle. Generally, the initial and final points of such a curve will lie in the same fiber, and be related by an element of G (the symmetry group acting on F). This element is the *holonomy* of C . Because it actually lies in the group (as opposed to the fiber), it is independent of the choice of gauge. We can obtain an explicit expression for the holonomy in terms of the connection A : an infinitesimal motion v in M gets lifted to $v + A_v$ in E , so the initial fiber coordinate gets multiplied by the group element $1 + A_v$. The total shift in the fiber when we go around a closed curve is

$$W = P \exp \oint_C A \cdot dx, \quad (3.A.12)$$

an object known in gauge theory as the Wilson line integral. This expression should be compared to Eq. (3.A.9); there are two differences that we should explain. The first is that Eq. (3.A.12) has no i in the exponent; this is because we have taken A to be Lie algebra valued, and thus anti-Hermitian, so the i has been absorbed into A . The second difference is the P in Eq. (3.A.12), which tells us how to order the non-abelian matrices $A(x)$ in expanding the exponential. The path ordering expressed by P means that all matrices A in the expansion of the exponential are to be ordered as the points in the path are, with later times on the right. If $x(t)$ is a parameterization of the curve C , the explicit expansion is

$$\begin{aligned} W &= P \exp \int A \cdot \dot{x} dt \\ &= 1 + \int_0^1 dt A \cdot \dot{x} + \int_0^1 dt \int_0^t dt' A(t) \cdot \dot{x}(t) A(t') \cdot \dot{x}(t') + \dots \end{aligned} \quad (3.A.13)$$

When the gauge potential is Abelian, the holonomy of Eq. (3.A.12) is the same as Eq. (3.A.9); everything commutes, so no path ordering is necessary. Furthermore, we can use Stokes' theorem (which does not apply in the non-abelian case) to relate the holonomy of a curve C to the magnetic flux enclosed by the curve. For our old friend the magnetic monopole, the holonomy associated to a closed curve C on the sphere is equal to $\exp ig\Omega$, where Ω is the solid angle subtended by C .

To summarize, there are gauge structures that give perfectly sensible rules for parallel transport which cannot be specified by a gauge potential A that is everywhere regular. We may specify A 's on different patches, and rules (gauge transformations) to identify them on the overlaps. Fiber bundles provide the most natural context in which to view such gauge potentials, whose singularities are closely related to the topology of the associated bundle. As we have seen, magnetic monopoles of different charge g generate different $U(1)$ bundles over a sphere. On the other hand, one can show that any two gauge potentials having the same total magnetic flux out of the

sphere, lead to topologically equivalent fiber bundles. Thus, the number of Dirac quanta of magnetic charge fully characterizes the topological class of the associated $U(1)$ bundle. The topological classification of fiber bundles in general is the subject of the theory of characteristic classes⁴.

This completes our short introduction to fiber bundles. A complementary approach to some of the topics we have covered is contained in the article by Aitchison [7.1]. Readers desiring a more formal and detailed treatment may wish to consult Ref. 4.

- [1] S. Coleman, "The Magnetic Monopole Fifty Years Later," in A. Zichichi, ed., *The Unity of the Fundamental Interactions* (New York: Plenum Press, 1983) pp.21-117.
- [2] R. Bott and L. Tu, *Differential Forms in Algebraic Topology* (New York: Springer-Verlag, 1982).
- [3] T.T. Wu and C.N. Yang, *Nucl. Phys.* **B107** (1976) 365.
- [4] N. Steenrod, *The Topology of Fibre Bundles* (Princeton: University Press, 1951);
D. Husemoller, *Fibre Bundles* (New York: Springer-Verlag, 1975).

Quantal phase factors accompanying adiabatic changes

By M. V. BERRY, F.R.S.

*H. H. Wills Physics Laboratory, University of Bristol,
Tyndall Avenue, Bristol BS8 1TL, U.K.*

(Received 13 June 1983)

A quantal system in an eigenstate, slowly transported round a circuit C by varying parameters \mathbf{R} in its Hamiltonian $\hat{H}(\mathbf{R})$, will acquire a geometrical phase factor $\exp\{iy(C)\}$ in addition to the familiar dynamical phase factor. An explicit general formula for $y(C)$ is derived in terms of the spectrum and eigenstates of $\hat{H}(\mathbf{R})$ over a surface spanning C . If C lies near a degeneracy of \hat{H} , $y(C)$ takes a simple form which includes as a special case the sign change of eigenfunctions of real symmetric matrices round a degeneracy. As an illustration $y(C)$ is calculated for spinning particles in slowly-changing magnetic fields; although the sign reversal of spinors on rotation is a special case, the effect is predicted to occur for bosons as well as fermions, and a method for observing it is proposed. It is shown that the Aharonov–Bohm effect can be interpreted as a geometrical phase factor.

1. INTRODUCTION

Imagine a quantal system whose Hamiltonian \hat{H} describes the effects of an unchanging environment, and let the system be in a stationary state. If the environment, and hence \hat{H} , is slowly altered, it follows from the adiabatic theorem (Messiah 1962) that at any instant the system will be in an eigenstate of the instantaneous \hat{H} . In particular, if the Hamiltonian is returned to its original form the system will return to its original state, apart from a phase factor. This phase factor is observable by interference if the cycled system is recombined with another that was separated from it at an earlier time and whose Hamiltonian was kept constant.

My purpose here is to explain how the phase factor contains a circuit-dependent component $\exp(iy)$ in addition to the familiar dynamical component $\exp(-iEt/\hbar)$ which accompanies the evolution of any stationary state. A general formula for y in terms of the eigenstates of \hat{H} will be obtained in § 2. If the circuit is close to a degeneracy in the spectrum of \hat{H} , y takes a particularly simple form which will be derived in § 3; this contains, as a special case, the sign change around a degeneracy of the eigenstates of a system whose Hamiltonian is real as well as Hermitian (Herzberg & Longuet-Higgins 1963; Longuet-Higgins 1975; Mead 1979; Mead & Truhlar 1979; Mead 1980a,b; Berry & Wilkinson 1984).

A particle of any spin in an eigenstate of a slowly-rotated magnetic field is another case where y can be calculated explicitly (§ 4), and gives predictions that could be

tested experimentally. This phase factor exists for bosons as well as fermions. A special case is the sign change of spinors slowly rotated by 2π , predicted by Aharonov & Susskind (1967); this will be shown to be different from the dynamical phase factors measured in experiments on precessing neutrons (reviewed by Silverman 1980).

Finally, it is shown in § 5 that physical effects of magnetic vector potentials in the absence of fields, predicted by Aharonov & Bohm (1959) and observed by Chambers (1960), can be understood as special cases of the geometrical phase factor.

2. GENERAL FORMULA FOR PHASE FACTOR

Let the Hamiltonian \hat{H} be changed by varying parameters $\mathbf{R} = (X, Y, \dots)$ on which it depends. Then the excursion of the system between times $t = 0$ and $t = T$ can be pictured as transport round a closed path $\mathbf{R}(t)$ in parameter space, with Hamiltonian $\hat{H}(\mathbf{R}(t))$ and such that $\mathbf{R}(T) = \mathbf{R}(0)$. The path will henceforth be called a circuit and denoted by C. For the adiabatic approximation to apply, T must be large.

The state $|\psi(t)\rangle$ of the system evolves according to Schrödinger's equation

$$\hat{H}(\mathbf{R}(t)) |\psi(t)\rangle = i\hbar |\dot{\psi}(t)\rangle. \quad (1)$$

At any instant, the natural basis consists of the eigenstates $|n(\mathbf{R})\rangle$ (assumed discrete) of $\hat{H}(\mathbf{R})$ for $\mathbf{R} = \mathbf{R}(t)$, that satisfy

$$\hat{H}(\mathbf{R}) |n(\mathbf{R})\rangle = E_n(\mathbf{R}) |n(\mathbf{R})\rangle, \quad (2)$$

with energies $E_n(\mathbf{R})$. This eigenvalue equation implies no relation between the phases of the eigenstates $|n(\mathbf{R})\rangle$ at different \mathbf{R} . For present purposes any (differentiable) choice of phases can be made, provided $|n(\mathbf{R})\rangle$ is single-valued in a parameter domain that includes the circuit C.

Adiabatically, a system prepared in one of these states $|n(\mathbf{R}(0))\rangle$ will evolve with \hat{H} and so be in the state $|n(\mathbf{R}(t))\rangle$ at t .

Thus $|\psi\rangle$ can be written as

$$|\psi(t)\rangle = \exp \left\{ \frac{-i}{\hbar} \int_0^t dt' E_n(\mathbf{R}(t')) \right\} \exp(i\gamma_n(t)) |n(\mathbf{R}(t))\rangle. \quad (3)$$

The first exponential is the familiar dynamical phase factor. In this paper the object of attention is the second exponential. The crucial point will be that its phase $\gamma_n(t)$ is *non-integrable*; γ_n cannot be written as a function of \mathbf{R} and in particular is not single-valued under continuation around a circuit, i.e. $\gamma_n(T) \neq \gamma_n(0)$.

The function $\gamma_n(t)$ is determined by the requirement that $|\psi(t)\rangle$ satisfy Schrödinger's equation, and direct substitution of (3) into (1) leads to

$$\dot{\gamma}_n(t) = i \langle n(\mathbf{R}(t)) | \nabla_{\mathbf{R}} n(\mathbf{R}(t)) \rangle \cdot \dot{\mathbf{R}}(t). \quad (4)$$

Phase factors accompanying adiabatic changes

47

The total phase change of $|\psi\rangle$ round C is given by

$$|\psi(T)\rangle = \exp(i\gamma_n(C)) \exp\left\{\frac{-i}{\hbar} \int_0^T dt E_n(\mathbf{R}(t))\right\} |\psi(0)\rangle, \quad (5)$$

where the *geometrical phase change* is

$$\gamma_n(C) = i \oint_C \langle n(\mathbf{R}) | \nabla_{\mathbf{R}} n(\mathbf{R}) \rangle \cdot d\mathbf{R}. \quad (6)$$

Thus $\gamma_n(C)$ is given by a circuit integral in parameter space and is independent of how the circuit is traversed (provided of course that this is slow enough for the adiabatic approximation to hold). The normalization of $|n\rangle$ implies that $\langle n | \nabla_{\mathbf{R}} n \rangle$ is imaginary, which guarantees that γ_n is real.

Direct evaluation of $|\nabla_{\mathbf{R}} n\rangle$ requires a locally single-valued basis for $|n\rangle$ and can be awkward. Such difficulties are avoided by transforming the circuit integral (6) into a surface integral over any surface in parameter space whose boundary is C. In order to employ familiar vector calculus, parameter space will be considered as three-dimensional, and this will turn out to be the important case in applications; the generalization to higher dimensions will be outlined at the end of this section.

Stokes's theorem applied to (6) gives, in an obvious abbreviated notation.

$$\gamma_n(C) = -\text{Im} \iint_C d\mathbf{S} \cdot \nabla \times \langle n | \nabla n \rangle, \quad (7a)$$

$$= -\text{Im} \iint_C d\mathbf{S} \cdot (\nabla n | \times | \nabla n \rangle, \quad (7b)$$

$$= -\text{Im} \iint_C d\mathbf{S} \cdot \sum_{m \neq n} \langle \nabla n | m \rangle \times \langle m | \nabla n \rangle, \quad (7c)$$

where $d\mathbf{S}$ denotes area element in \mathbf{R} space and the exclusion in the summation is justified by $\langle n | \nabla n \rangle$ being imaginary. The off-diagonal elements are obtained from (2) as

$$\langle m | \nabla n \rangle = \langle m | \nabla \hat{H} | n \rangle / (E_n - E_m), \quad m \neq n. \quad (8)$$

Thus γ_n can be expressed as

$$\gamma_n(C) = - \iint_O d\mathbf{S} \cdot V_n(\mathbf{R}), \quad (9)$$

where

$$V_n(\mathbf{R}) \equiv \text{Im} \sum_{m \neq n} \frac{\langle n(\mathbf{R}) | \nabla_{\mathbf{R}} \hat{H}(\mathbf{R}) | m(\mathbf{R}) \rangle \times \langle m(\mathbf{R}) | \nabla_{\mathbf{R}} \hat{H}(\mathbf{R}) | n(\mathbf{R}) \rangle}{(E_m(\mathbf{R}) - E_n(\mathbf{R}))^2}. \quad (10)$$

Obviously $\gamma_n(C)$ is zero for a circuit which retraces itself and so encloses no area.

Equations (9) and (10) embody the central results of this paper. Because the dependence on $|\nabla n\rangle$ has been eliminated, phase relations between eigenstates with different parameters are now immaterial, and (as is evident from the form of (10)), it is no longer necessary to choose $|m\rangle$ and $|n\rangle$ to be single-valued in \mathbf{R} : any solutions of (2) may be employed without affecting the value of V_n . This is a surprising conclusion, as can be seen by comparing (9) with (7a) which show that V_n is the curl of a vector, $\langle n | \nabla n \rangle$, and $\langle n | \nabla n \rangle$ certainly does depend on the choice of phase

of the (single-valued) eigenstate $|n(\mathbf{R})\rangle$. The dependence on phase is of the following kind: if $|n\rangle \rightarrow \exp\{i\mu(\mathbf{R})\}|n\rangle$ then $\langle n|\nabla n\rangle \rightarrow \langle n|\nabla n\rangle + i\nabla\mu$ (in another context the importance of such gauge transformations has been emphasized by Wu & Yang (1975)). Thus the vector is not unique but its curl is. The quantity \mathbf{V}_n is analogous to a ‘magnetic field’ (in parameter space) whose ‘vector potential’ is $\text{Im}\langle n|\nabla n\rangle$. In Appendix A it is shown directly from (10) that $\nabla \cdot \mathbf{V}_n$ vanishes, thus confirming that (9) gives a unique value for $\gamma_n(\mathcal{C})$.

Using perturbation theory, Mead & Truhlar (1979) obtained essentially the formulae (9) and (10) for an infinitesimal circuit, in a study of molecular electronic states which (in the Born–Oppenheimer approximation) depend parametrically on nuclear coordinates. Their phase factor was not intended to apply to a $|\psi\rangle$ that evolves slowly under the time-dependent Schrödinger equation, but to the variation of eigenstates $|n\rangle$ under a particular phase-continuation rule in \mathbf{R} -space which can be shown to give the same result.

In parameter spaces of higher dimension, Stokes’s theorem cannot be employed to transform the circuit integral (6). The appropriate generalization, provided by the theory of differential forms (see, for example, Arnold 1978, chap. 7), transforms (6) into the integral of a 2-form over a surface bounded by \mathcal{C} . The surprising result (10) can now be expressed as follows: independently of the choice of phases of the eigenstates, there exists in parameter space a *phase 2-form*, which gives $\gamma(\mathcal{C})$ when integrated over any surface spanning \mathcal{C} . This 2-form is obtained from (10) by replacing ∇ by the exterior derivative d and \times by the wedge product \wedge . The validity of this generalization is consistent with the observation that in the three-dimensional version there are infinitely many choices of interpolating Hamiltonian (and hence of parameter spaces) on the surfaces bounded by \mathcal{C} , and the geometrical phase factor is independent of the choice.

Professor Barry Simon (1983), commenting on the original version of this paper, points out that the geometrical phase factor has a mathematical interpretation in terms of holonomy, with the phase two-form emerging naturally (in the form (7b)) as the curvature (first Chern class) of a Hermitian line bundle.

3. DEGENERACIES

The energy denominators in (10) show that if the circuit \mathcal{C} lies close to a point \mathbf{R}^* in parameter space at which the state n is involved in a degeneracy, then $\mathbf{V}_n(\mathbf{R})$, and hence $\gamma_n(\mathcal{C})$, is dominated by the terms m corresponding to the other states involved. We shall consider the commonest situation, where the degeneracy involves only two states, to be denoted $+$ and $-$, where $E_+(\mathbf{R}) \geq E_-(\mathbf{R})$. For \mathbf{R} near \mathbf{R}^* , $\hat{H}(\mathbf{R})$ can be expanded to first order in $\mathbf{R} - \mathbf{R}^*$, and

$$\mathbf{V}_+(\mathbf{R}) = \text{Im} \frac{\langle +(\mathbf{R}) | \nabla \hat{H}(\mathbf{R}^*) | -(\mathbf{R}) \rangle \times \langle -(\mathbf{R}) | \nabla \hat{H}(\mathbf{R}^*) | +(\mathbf{R}) \rangle}{(E_+(\mathbf{R}) - E_-(\mathbf{R}))^2}. \quad (11)$$

Obviously $\mathbf{V}_-(\mathbf{R}) = -\mathbf{V}_+(\mathbf{R})$, so that $\gamma_-(\mathcal{C}) = -\gamma_+(\mathcal{C})$.

Phase factors accompanying adiabatic changes

49

Without essential loss of generality we can take $E_{\pm}(\mathbf{R}^*) = 0$ and $\mathbf{R}^* = 0$. $H(\mathbf{R})$ can be represented by a 2×2 Hermitian matrix coupling the two states. The most general such matrix satisfying the given conditions depends on three parameters X, Y, Z which will be taken as components of \mathbf{R} , and by linear transformation in \mathbf{R} -space can be brought into the following standard form

$$\hat{H}(\mathbf{R}) = \frac{1}{2} \begin{bmatrix} Z & X - iY \\ X + iY & -Z \end{bmatrix}. \quad (12)$$

The eigenvalues are

$$E_+(\mathbf{R}) = -E_-(\mathbf{R}) = \frac{1}{2}(X^2 + Y^2 + Z^2)^{\frac{1}{2}} = \frac{1}{2}R. \quad (13)$$

Thus the degeneracy is an isolated point at which all three parameters vanish. This illustrates an old result of Von Neumann & Wigner (1929): for generic Hamiltonians (Hermitian matrices), it is necessary to vary three parameters in order to make a degeneracy occur accidentally, that is, not on account of symmetry. Alternatively stated, degeneracies have co-dimension three.

The form (12) was chosen to exploit the fact that

$$\nabla \hat{H} = \frac{1}{2}\hat{\sigma}, \quad (14)$$

where the components $\hat{\sigma}_X, \hat{\sigma}_Y, \hat{\sigma}_Z$ of the vector operator $\hat{\sigma}$ are the Pauli spin matrices. When evaluating the matrix elements in (11) it greatly simplifies the calculations to take advantage of the isotropy of spin and temporarily rotate axes so that the Z -axis points along \mathbf{R} , and to employ the following relations, which come from the commutation laws between the components of $\hat{\sigma}$:

$$\hat{\sigma}_X |\pm\rangle = |\mp\rangle, \quad \hat{\sigma}_Y |\pm\rangle = \pm i|\mp\rangle, \quad \hat{\sigma}_Z |\pm\rangle = \pm |\pm\rangle. \quad (15)$$

With these rotated axes, (11) gives

$$\left. \begin{aligned} V_{X+} &= (\text{Im} \langle + | \hat{\sigma}_Y | - \rangle \langle - | \hat{\sigma}_Z | + \rangle) / 2R^2 = 0, \\ V_{Y+} &= (\text{Im} \langle + | \hat{\sigma}_Z | - \rangle \langle - | \hat{\sigma}_X | + \rangle) / 2R^2 = 0, \\ V_{Z+} &= \text{Im} \langle + | \hat{\sigma}_X | - \rangle \langle - | \hat{\sigma}_Y | + \rangle = 1/2R^2. \end{aligned} \right\} \quad (16)$$

Reverting to unrotated axes, we obtain

$$V_+(\mathbf{R}) = R/2R^3. \quad (17)$$

Now use of (9) shows that the phase change $\gamma_+(C)$ is the flux through C of the magnetic field of a monopole with strength $-\frac{1}{2}$ located at the degeneracy. Thus we obtain the pleasant result, valid for the natural choice (12) of standard form for \hat{H} , that the geometrical phase factor associated with C is

$$\exp \{i\gamma_{\pm}(C)\} = \exp \{ \mp \frac{1}{2}i\Omega(C) \}, \quad (18)$$

where $\Omega(C)$ is the *solid angle* that C subtends at the degeneracy; Ω is, in a sense, a measure of the *view* of the circuit as seen from the degeneracy. The phase factor is

independent of the choice of surface spanning C , because Ω can change only in multiples of 4π (when the surface is deformed to pass through the degeneracy).

An important special case of (18) occurs when C consists entirely of *real* Hamiltonians and so is confined to the plane $Y = 0$ (cf. (12)). The energy levels E_{\pm} intersect conically in the space E, X, Z , whose origin, where the degeneracy occurs, is a ‘diabolical point’ of the type recently studied by Berry & Wilkinson (1984) in the spectra of triangles. This illustrates the result that for real symmetric matrices, degeneracies have co-dimension two: see Appendix 10 of Arnold 1978. If C encloses the degeneracy, $\Omega = \pm 2\pi$; if not, $\Omega = 0$. Thus the phase factor (18) is

$$\begin{aligned}\exp\{i\gamma_{\pm}(C)\} &= -1, && \text{if } C \text{ encircles the degeneracy,} \\ &= +1, && \text{otherwise,}\end{aligned}\tag{19}$$

which expresses the sign changes of real wavefunctions as a degeneracy involving them is encircled, a phenomenon first described by Herzberg & Longuet-Higgins (1963). (Sign changes are not restricted to circuits involving real Hamiltonians: (18) shows that the phase factor is -1 if C lies in *any* plane through the degeneracy and encircles it.)

Confirmation of the correctness of (17) can be obtained without the rotation-of-axes trick, by a lengthy calculation of (11) involving explicit formulae for the eigenvectors $|\pm(\mathbf{R})\rangle$ of the matrix (12). Alternatively, direct continuation of the eigenvectors may be attempted. This cannot be accomplished for all circuits by means of (6) because it is not possible to construct eigenvectors that are globally single-valued continuous functions of \mathbf{R} ; multivaluedness can be reduced to singular lines connecting the degeneracy with infinity, and in the analogue $V(\mathbf{R})$ these appear as Dirac strings attached to the monopole. Such approaches obscure the simplicity and essential isotropy of the solid-angle result (17).

Using topological arguments not involving explicit formulae for $\gamma_n(C)$, Stone (1976) proved that if C is expanded from one point \mathbf{R} and contracted on to another so as to sweep out a surface enclosing a degeneracy, then the geometrical phase factor traverses a circle in its Argand plane. This property (which follows easily from (18)), is the Hermitian generalization of the sign-reversal test for degeneracy.

4. SPINS IN MAGNETIC FIELDS

A particle with spin s (integer or half-integer) interacts with a magnetic field \mathbf{B} via the Hamiltonian

$$\hat{H}(\mathbf{B}) = \kappa\hbar\mathbf{B}\cdot\hat{\mathbf{s}},\tag{20}$$

where κ is a constant involving the gyromagnetic ratio and $\hat{\mathbf{s}}$ is the vector spin operator with $2s+1$ eigenvalues n with integer spacing and that lie between $-s$ and $+s$. The eigenvalues are

$$E_n(\mathbf{B}) = \kappa\hbar B n,\tag{21}$$

and so there is a $(2s+1)$ -fold degeneracy when $\mathbf{B} = 0$. (The special case $s = \frac{1}{2}$ reproduces the two-fold degeneracy considered in the last section.) We consider

the components of \mathbf{B} as the parameters \mathbf{R} in our previous analysis, and calculate the phase change $\gamma_n(C)$ of an eigenstate $|n, s(\mathbf{B})\rangle$ of $\hat{\mathbf{s}}$ in the direction along \mathbf{B} , as \mathbf{B} is slowly varied (and hence the spin rotated) round a circuit C .

The vector $V_n(\mathbf{B})$ as given by (10) can be expressed by using (20) and (21) as

$$V_n(\mathbf{B}) = \frac{\text{Im}}{B^2} \sum_{m \neq n} \frac{\langle n, s(\mathbf{B}) | \hat{s} | m, s(\mathbf{B}) \rangle \times \langle m, s(\mathbf{B}) | \hat{s} | n, s(\mathbf{B}) \rangle}{(m - n)^2}. \quad (22)$$

To evaluate the matrix elements we again temporarily rotate axes so that the Z -axis points along \mathbf{B} , and employ the following generalization of (15):

$$\left. \begin{aligned} (\hat{s}_X + i\hat{s}_Y) |n, s\rangle &= [s(s+1) - n(n+1)]^{\frac{1}{2}} |n+1, s\rangle, \\ (\hat{s}_X - i\hat{s}_Y) |n, s\rangle &= [s(s+1) - n(n-1)]^{\frac{1}{2}} |n-1, s\rangle, \\ s_Z |n, s\rangle &= n |n, s\rangle. \end{aligned} \right\} \quad (23)$$

It is clear that only states with $m = n \pm 1$ are coupled with $|n\rangle$ in (22), and that V_x and V_y are zero because they involve off-diagonal elements of \hat{s}_Z . To find V_z , we make use of (23) to obtain

$$\left. \begin{aligned} \langle n \pm 1, s | s_X | n, s \rangle &= \frac{1}{2}[s(s+1) - n(n \pm 1)]^{\frac{1}{2}}, \\ \langle n \pm 1, s | s_Y | n, s \rangle &= \mp \frac{1}{2}i[s(s+1) - n(n \pm 1)]^{\frac{1}{2}}, \end{aligned} \right\} \quad (24)$$

then (22) gives

$$\begin{aligned} V_{Zn} &= \frac{\text{Im}}{B^2} \{ \langle n | s_X | n+1 \rangle \langle n+1 | s_Y | n \rangle - \langle n | s_Y | n+1 \rangle \langle n+1 | s_X | n \rangle \\ &\quad + \langle n | s_X | n-1 \rangle \langle n-1 | s_Y | n \rangle - \langle n | s_Y | n-1 \rangle \langle n-1 | s_X | n \rangle \} \\ &= \frac{n}{B^2}. \end{aligned} \quad (25)$$

Reverting to unrotated axes, we obtain

$$V_n(\mathbf{B}) = n\mathbf{B}/B^3. \quad (26)$$

Now, use of (9) shows that $\gamma_n(C)$ is the flux through C of the ‘magnetic field’ of a monopole $-n$ located at the origin of magnetic field space. Thus the geometrical phase factor is

$$\exp\{i\gamma_n(C)\} = \exp\{-in\Omega(C)\}, \quad (27)$$

where $\Omega(C)$ is the solid angle that C subtends at $\mathbf{B} = 0$. Note that γ_n depends only on the eigenvalue n of the spin component along \mathbf{B} and not on the spin s of the particle, so that γ_n is insensitive to the strength $2s+1$ of the degeneracy at $\mathbf{B} = 0$.

It follows from (27) that any phase change can be produced by varying \mathbf{B} round a suitable circuit. For fermions (half-integer n), a whole turn of \mathbf{B} (rotation through 2π in a plane, giving $\Omega = 2\pi$) produces a phase factor -1 . In the special case $n = \frac{1}{2}$ this shows that the sign change of spinors on rotation and the sign change of wavefunctions round a degeneracy have the same mathematical origin. For bosons (integer n), a whole turn of \mathbf{B} produces a phase factor $+1$. To produce a sign change,

different circuits are required; if $n = 1$, for example, varying \mathbf{B} round a cone of semiangle 60° will give $\Omega = \gamma = \pi$ and hence a phase factor -1 .

The following experiment could be carried out to test the predictions embodied in (27). A polarized monoenergetic beam of particles in spin state n along a magnetic field \mathbf{B} is split into two. Along the path of one beam \mathbf{B} is kept constant. Along the path of the other beam, \mathbf{B} is kept constant in magnitude but its direction is varied slowly (in comparison with the dynamical precession frequency) round a circuit C subtending a solid angle Ω , the field being generated by an arrangement enabling Ω to be changed. The beams are then combined and the count rate at a detector is measured as a function of Ω . The dynamical phase factor (the second exponential in (5) is the same for both beams because the energy $E_n(\mathbf{B})$ (21) is insensitive to the direction of \mathbf{B} . There will in addition be a propagation phase factor which can be made unity by adjusting the path-length of one of the beams when $\Omega = 0$. The resulting fringes occur as a consequence of the geometrical phase factor. If C is a circuit round a cone of semiangle θ , the predicted intensity contrast is

$$I(\theta) = \cos^2(n\pi(1 - \cos\theta)). \quad (28)$$

I wish to emphasize that this proposed experiment is different from those carried out by Rauch *et al.* (1975, 1978) and Werner *et al.* (1975) (see Silverman 1980) with *unpolarized* neutrons in a *constant* magnetic field. Those neutrons were not in an eigenstate, and their phase changed dynamically, rather than geometrically, under the Hamiltonian (20) (with \mathbf{B} along Z and $\hat{\sigma}$ replacing \hat{s}) according to the evolution operator

$$\exp(-i\hat{H}t/\hbar) = \exp(-Bkt\hat{\sigma}_Z) = \cos \frac{1}{2}\kappa Bt \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + i \sin \frac{1}{2}\kappa Bt \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}. \quad (29)$$

The sign changed whenever $\frac{1}{2}\kappa Bt$ was an odd multiple of π , and this was interpreted on the basis of precession theory as corresponding to odd numbers of complete rotations about \mathbf{B} .

5. AHARONOV-BOHM EFFECT

Consider a magnetic field consisting of a single line with flux Φ . For positions \mathbf{R} not on the flux line, the field is zero but there must be a vector potential $\mathbf{A}(\mathbf{R})$ satisfying

$$\oint_C \mathbf{A}(\mathbf{R}) \cdot d\mathbf{R} = \Phi, \quad (30)$$

for circuits C threaded by the flux line. Aharonov & Bohm (1959) showed that in quantum mechanics such vector potentials have physical significance even though they correspond to zero field. I shall now show how their effect can be interpreted as a geometrical phase change of the type described in § 2.

Let the quantal system consist of particles with charge q confined to a box situated at \mathbf{R} and not penetrated by the flux line (figure 1). In the absence of flux

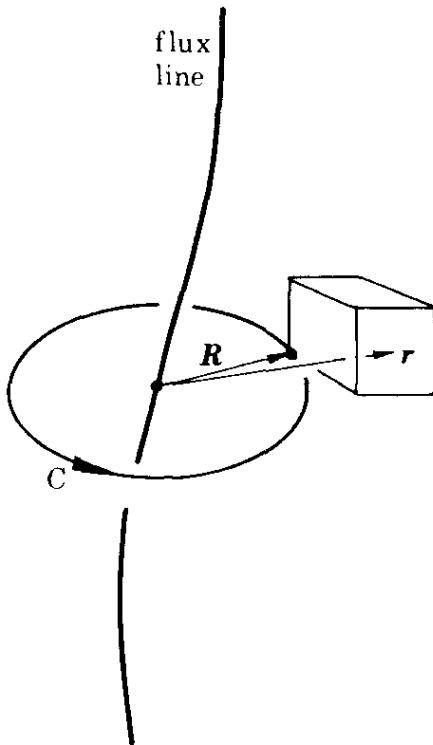


FIGURE 1. Aharonov-Bohm effect in a box transported round a flux line.

($A = 0$), the particle Hamiltonian depends on position \hat{r} and conjugate momentum \hat{p} as follows:

$$\hat{H} = H(\hat{p}, \hat{r} - \mathbf{R}), \quad (31)$$

and the wavefunctions have the form $\psi_n(\mathbf{r} - \mathbf{R})$ with energies E_n independent of \mathbf{R} . With non-zero flux, the states $|n(\mathbf{R})\rangle$ satisfy

$$H(\hat{p} - q\mathbf{A}(\hat{r}), \hat{r} - \mathbf{R})|n(\mathbf{R})\rangle = E_n |n(\mathbf{R})\rangle, \quad (32)$$

an equation whose exact solutions are obtained by multiplying ψ_n by an appropriate Dirac phase factor, giving

$$\langle \mathbf{r} | n(\mathbf{R}) \rangle = \exp \left\{ \frac{iq}{\hbar} \int_{\mathbf{R}}^{\mathbf{r}} d\mathbf{r}' \cdot \mathbf{A}(\mathbf{r}') \right\} \psi_n(\mathbf{r} - \mathbf{R}). \quad (33)$$

These solutions are single-valued in \mathbf{r} and (locally) in \mathbf{R} . The energies are unaffected the vector potential.

Now let the box be transported round a circuit C threaded by the flux line; in this particular case it is not necessary that the transport be adiabatic. After completion of the circuit there will be a geometrical phase change that can be calculated from (6) and (33) by using

$$\begin{aligned} \langle n(\mathbf{R}) | \nabla_{\mathbf{R}} n(\mathbf{R}) \rangle &= \iiint d^3 \mathbf{r} \psi_n^*(\mathbf{r} - \mathbf{R}) \left\{ \frac{-iq}{\hbar} \mathbf{A}(\mathbf{R}) \psi_n(\mathbf{r} - \mathbf{R}) + \nabla_{\mathbf{R}} \psi_n(\mathbf{r} - \mathbf{R}) \right\} \\ &= -iq \mathbf{A}(\mathbf{R}) / \hbar. \end{aligned} \quad (34)$$

(The vanishing of the second term in braces follows from the normalization of ψ_n .) Evidently in this example the analogy between $\text{Im} \langle n | \nabla n \rangle$ and a magnetic vector potential becomes a reality. Thus

$$\gamma_n(C) = \frac{q}{\hbar} \oint_C \mathbf{A}(\dot{\mathbf{R}}) \cdot d\mathbf{R} = q\Phi/\hbar, \quad (35)$$

which shows that the phase factor is independent of n , and also of C if this winds once round the flux line. The phase factor could be observed by interference between the particles in the transported box and those in a box which was not transported round the circuit.

In elementary presentations of the Aharonov–Bohm effect (including its anticipation by Ehrenburg & Siday 1949), the Dirac phase factor is often invoked in comparing systems passing opposite sides of a flux line. Such invocations are subject to the objection that the wavefunction thus obtained is not single-valued. One way to avoid this objection is by summation over all contributions (whirling waves) representing different windings round the flux line (Schulman 1981; Berry 1980; Morandi & Menossi 1984). Another way, adopted in the original paper by Aharonov & Bohm, is to solve Schrödinger’s equation exactly for scattering in the flux line’s vector potential. The argument of the preceding paragraphs, which employs the geometrical phase factor, is a third way of obtaining the Aharonov–Bohm effect by using only single-valued wavefunctions.

Mead (1980a, b), employs the term ‘molecular Aharonov–Bohm effect’ in a different context, to describe how degeneracies in electron energy levels affect the spectrum of nuclear vibrations. He explains two options, both leading to the same vibration spectrum. The first option is to continue the electronic states round degeneracies (in the space of nuclear coordinates) in the manner described in this paper, thus causing the electronic wavefunctions to be multi-valued, with a compensating multi-valuedness in the nuclear states, which must be incorporated into their boundary conditions. The alternative is to enforce single-valuedness on the electronic (and hence also the nuclear) states, and this introduces a vector potential into the Schrödinger equation for nuclear motion. In general one may expect such effects whenever an isolated system is considered as being divided into two interacting parts, each slaved to a different aspect of the other (in the molecular case, electron states are slaved to nuclear coordinates, and nuclear states are slaved to the electronic states and wavefunctions). The systems considered in this paper might be regarded as a special case, in which the coupling is with ‘the rest of the Universe’ (including us as observers). The only role of the rest of the Universe is to provide a Hamiltonian with slowly-varying parameters, thus forcing the system to evolve adiabatically with phase continuation governed by the time-dependent Schrödinger equation.

6. DISCUSSION

A system slowly transported round a circuit will return in its original state; this is the content of the adiabatic theorem. Moreover its internal clocks will register the passage of time; this can be regarded as the meaning of the dynamical phase factor. The remarkable and rather mysterious result of this paper is that in addition the system records its history in a deeply geometrical way, whose natural formulation (9) and (10) involves phase functions hidden in parameter-space regions which the system has not visited.

The total phase of the transported state (5) is dominated by the dynamical part, because $T \rightarrow \infty$ in the adiabatic limit, and it might be thought that this must overwhelm the geometrical phase γ_n and make its physical effects difficult to detect. This objection can be met by observing that the strengths of non-adiabatic transitions are exponentially small in T if \hat{H} changes smoothly (Hwang & Pechukas 1977), so that essentially adiabatic evolution can occur even when the dynamical phase is only a few times greater than 2π .

As we saw in § 3, degeneracies in the spectrum of $\hat{H}(\mathbf{R})$ are the singularities of the vector $\mathbf{V}(\mathbf{R})$ (equation (10)) in parameter space, and so have an important effect on the geometrical phase factor. This is reminiscent of the part played by singularities of an analytic function, but the analogy is imperfect: if $\gamma(C)$ were completely singularity-determined, $\mathbf{V}(\mathbf{R})$ would be the sum of the ‘magnetic fields’ of ‘monopoles’ situated at the degeneracies (cf. (17)) and so would have zero curl, which is not the case (zero curl, unlike zero divergence, is not a property which is invariant under deformations of \mathbf{R} space, and in the general case the sources of \mathbf{V} are not just monopoles but also ‘currents’ distributed continuously in parameter space). A closer analogy is with wavefront dislocation lines, which are phase singularities of complex wavefunctions in three-dimensional position space (Nye & Berry 1974; Nye 1981; Berry 1981), that dominate the geometry of wavefronts without completely determining them.

In view of the emphasis on degeneracies as organizing centres for phase changes, it is worth remarking that close approach of energy levels is not a necessary condition for the existence of nontrivial geometrical phase factors. Indeed, our examples have shown that $\gamma(C)$ can be non-zero even if C involves isospectral deformations of $\hat{H}(\mathbf{R})$ (in the Aharonov–Bohm illustration, the levels E_n do not depend on \mathbf{R} at all).

The results obtained here are not restricted to quantum mechanics, but apply more generally, to the phase of eigenvectors of any Hermitian matrices under a natural continuation in parameter space. Therefore they have implications throughout wave physics. For example, the electromagnetic field of a single mode travelling along an optical fibre will change sign if the cross section of the fibre is slowly altered so that its path (in the space of shapes) surrounds a shape for which the spectrum of the Helmholtz equation is degenerate (such as one of the diabolical triangles discovered by Berry & Wilkinson 1984).

I thank Dr J. H. Hannay, Dr E. J. Heller and Dr B. R. Pollard for several suggestions, and Professor Barry Simon for showing me, before publication, his paper which comments on this one. This work was not supported by any military agency.

REFERENCES

- Aharonov, Y. & Bohm, D. 1959 *Phys. Rev.* **115**, 485–491.
 Aharonov, Y. & Susskind, L. 1967 *Phys. Rev.* **158**, 1237–1238.
 Arnold, V. I. 1978 *Mathematical methods of classical dynamics*. New York: Springer.
 Berry, M. V. 1980 *Eur. J. Phys.* **1**, 240–244.
 Berry, M. V. 1981 Singularities in waves and rays. In *Les Houches Lecture Notes for session XXXV* (ed. R. Balian, M. Kléman & J.-P. Poirier), pp. 453–543. Amsterdam: North-Holland.
 Berry, M. V. & Wilkinson, M. 1984 *Proc. R. Soc. Lond. A* **392**, 15–43.
 Chambers, R. G. 1960 *Phys. Rev. Lett.* **5**, 3–5.
 Ehrenburg, W. & Siday, R. E. 1949 *Proc. phys. Soc. B* **62**, 8–21.
 Herzberg, G. & Longuet-Higgins, H. C. 1963 *Discuss. Faraday Soc.* **35**, 77–82.
 Hwang, J.-T. & Pechukas, P. 1977 *J. chem. Phys.* **67**, 4640–4653.
 Longuet-Higgins, H. C. 1975 *Proc. R. Soc. Lond. A* **344**, 147–156.
 Mead, C. A. 1979 *J. chem. Phys.* **70**, 2276–2283.
 Mead, C. A. 1980a *Chem. Phys.* **49**, 23–32.
 Mead, C. A. 1980b *Chem. Phys.* **49**, 33–38.
 Mead, C. A. & Truhlar, D. G. 1979 *J. chem. Phys.* **70**, 2284–2296.
 Messiah, A. 1962 *Quantum mechanics*, vol. 2. Amsterdam: North-Holland.
 Morandi, G. & Menossi, E. 1984 *Nuovo Cim. B* (Submitted.)
 Nye, J. F. 1981 *Proc. R. Soc. Lond. A* **378**, 219–239.
 Nye, J. F. & Berry, M. V. 1974 *Proc. R. Soc. Lond. A* **336**, 165–190.
 Rauch, H., Wilfing, A., Bauspiess, W. & Bonse, U. 1978 *Z. Phys. B* **29**, 281–284.
 Rauch, H., Zeilinger, A., Badurek, G., Wilfing, A., Bauspiess, W. & Bonse, U. 1975 *Physics Lett. A* **54**, 425–427.
 Schulman, L. S. 1981 Techniques and Applications of Path Integration. New York: John Wiley.
 Silverman, M. P. 1980 *Eur. J. Phys.* **1**, 116–123.
 Simon, B. 1983 *Phys. Rev. Lett.* (In the press.)
 Stone, A. J. 1976 *Proc. R. Soc. Lond. A* **351**, 141–150.
 Von Neumann, J. & Wigner, E. P. 1929 *Phys. Z.* **30**, 467–470.
 Werner, S. A., Colella, R., Overhauser, A. W. & Eagen, C. F. 1975 *Phys. Rev. Lett.* **35**, 1053–1055.
 Wu, T. T. & Yang, C. N. 1975 *Phys. Rev. D* **12**, 3845–3857.

APPENDIX A

To show that $\gamma(C)$ is independent of the surface spanning C , it is necessary to prove that $V(R)$ (equation (10)) has zero divergence. This can be accomplished by expressing V in terms of the vector Hermitian operator \hat{B} defined by

$$\hat{B} \equiv -i \sum_n |\nabla n\rangle \langle n|. \quad (\text{A } 1)$$

From (8), the off-diagonal elements of \hat{B} are

$$\langle n | \hat{B} | m \rangle = -i \langle m | \nabla H | n \rangle / (E_n - E_m), \quad m \neq n. \quad (\text{A } 2)$$

Thus (10) becomes

$$V = \text{Im} \langle n | \hat{\mathbf{B}} \times \hat{\mathbf{B}} | n \rangle. \quad (\text{A } 3)$$

Now we can calculate the divergence:

$$\nabla \cdot V = \text{Im} \{ \langle \nabla n | \cdot \hat{\mathbf{B}} \times \hat{\mathbf{B}} | n \rangle + \langle n | \mathbf{B} \times \mathbf{B} \cdot | \nabla n \rangle + \langle n | \nabla \cdot (\hat{\mathbf{B}} \times \hat{\mathbf{B}}) | n \rangle \}, \quad (\text{A } 4)$$

Use of a consequence of (A 1), namely

$$|\nabla n\rangle = i\hat{\mathbf{B}}|n\rangle \quad (\text{A } 5)$$

gives

$$\nabla \cdot V = n(-\hat{\mathbf{B}} \cdot \hat{\mathbf{B}} \times \hat{\mathbf{B}} + \hat{\mathbf{B}} \times \hat{\mathbf{B}} \cdot \hat{\mathbf{B}})|n\rangle + \text{Im} \langle n | (\nabla \times \hat{\mathbf{B}} \cdot \hat{\mathbf{B}} - \hat{\mathbf{B}} \cdot \nabla \times \mathbf{B}) | n \rangle. \quad (\text{A } 6)$$

For the curl of \mathbf{B} , (A 1) and (A 5) give

$$\nabla \times \hat{\mathbf{B}} = +i \sum_n |\nabla n\rangle \times \langle \nabla n| = i \sum_n \hat{\mathbf{B}}|n\rangle \times \langle n|\hat{\mathbf{B}} = i\hat{\mathbf{B}} \times \hat{\mathbf{B}}, \quad (\text{A } 7)$$

whence $\nabla \cdot V$ vanishes by the dot-cross rule for triple products.

This result is valid everywhere except at the ‘monopole’ singularities arising from degeneracies.

Holonomy, the Quantum Adiabatic Theorem, and Berry's Phase

Barry Simon

Departments of Mathematics and Physics, California Institute of Technology, Pasadena, California 91125

(Received 18 October 1983)

It is shown that the "geometrical phase factor" recently found by Berry in his study of the quantum adiabatic theorem is precisely the holonomy in a Hermitian line bundle since the adiabatic theorem naturally defines a connection in such a bundle. This not only takes the mystery out of Berry's phase factor and provides calculational simple formulas, but makes a connection between Berry's work and that of Thouless *et al.* This connection allows the author to use Berry's ideas to interpret the integers of Thouless *et al.* in terms of eigenvalue degeneracies.

PACS numbers: 03.65.Db, 02.40.+m

Vector bundles and their integral invariants (Chern numbers) are already familiar to theoretical physicists because of their occurrence in classical Yang-Mills theories. Here I want to explain how they also enter naturally into non-relativistic quantum mechanics, especially in problems connected with condensed matter physics. If one has a Hermitian operator $\tilde{H}(x)$ depending smoothly on a parameter x , with an isolated nondegenerate eigenvalue $E(x)$ depending continuously on x , then $\{(x, \varphi) | \tilde{H}(x)\varphi = E(x)\varphi\}$ defines a line bundle over the parameter space. I will show that the twisting of this line bundle affects the phase of quantum mechanical wave functions.

Berry, in a beautiful recent paper,¹ discovered a striking phenomenon in the quantum adiabatic theorem.² That theorem says³ that if $H(t)$, $0 \leq t \leq 1$, is a family of Hermitian Hamiltonians, depending smoothly on t , and if $E(t)$ is a smooth function of t which is a simple eigenvalue of $H(t)$ isolated from the rest of the spectrum of $H(t)$ for each t , then the solution $\psi_T(s)$ of the time-dependent Schrödinger equation

$$i d\psi_T(s)/ds = H(s/T)\psi_T(s) \quad (1)$$

with $\psi_T(0) = \varphi_0$ where $H(0)\varphi_0 = E(0)\varphi_0$ has the property that as $T \rightarrow \infty$, $\psi_T(T)$ approaches the eigenvector φ_1 of $H(1)$ with $H(1)\varphi_1 = E(1)\varphi_1$. Berry asked the following question: Suppose that $\tilde{H}(x)$ is a multiple-parameter family and that $C(t)$ is a closed curve in parameter space, so that $\tilde{H}(C(t)) = H(t)$ obeys the hypotheses of the adiabatic theorem. Then that theorem says that φ_1 is just a phase factor times φ_0 and Berry asks, "What phase factor?" Surprisingly, the "obvious" guess

$$\varphi_1 = \exp[-i \int_0^T E(s/T) ds] \varphi_0$$

is wrong; rather, Berry finds

$$\varphi_1 = \exp[-i \int_0^T E(s/T) ds] \exp[i\gamma(C)] \varphi_0, \quad (2)$$

where $\gamma(C)$ is an extra phase which Berry extensively studies, and which he suggests could be experimentally measured.⁴

The purpose here is first to advertise what Berry calls a "remarkable and rather mysterious result," but more basically to try to take the mystery out of it by realizing that γ is an integral of a curvature so that Berry's phenomenon is essentially that of holonomy which is becoming quite familiar to theoretical physicists.⁵ This realization will allow us a more compact formula than that used principally by Berry, and one that is easier to compute with. Most importantly, it will give a close mathematical relationship between his work and that of Thouless *et al.*,⁶ so that Berry's interesting analysis of the relation of degeneracy to $\gamma(E)$ will allow us to interpret the TKN² integers in a new and interesting way.

To explain that γ is a holonomy, I begin by replacing $H(s)$ by $H(s) - E(s)$ which produces a trivial, computable phase change in the solution $\psi_T(s)$ of (1). Thus, without loss, we can take $E(s) = 0$. Define $\eta_T(s) = \psi_T(sT)$ so that η solves

$$i d\eta_T(s)/ds = TH(s)\eta_T(s), \quad (3)$$

and the adiabatic theorem says that $\eta_T(s)$ has a limit $\eta(s)$ with

$$H(s)\eta(s) = 0 \quad (4)$$

[since we have taken $E(s) = 0$]. If now $\tilde{H}(x)$ is a multiparameter family of Hermitian Hamiltonians, so that in some region M of parameter space, 0 is an isolated nondegenerate eigenvalue, then given any curve $C(t)$ and any choice η_0 of normalized zero-energy eigenvector of $H(C(0))$ (i.e., a choice of phase), the adiabatic limit yields a way of transporting η_0 along the curve $C(t)$, i.e., a connection. In this way, (2) is just an expression of the holonomy associated to this connection. So far this is just fancy words to de-

scribe Berry's discovery. However, there is a mathematically natural connection already long known in the situation of distinguished lines in Hilbert space. For given x , let X_x denote the zero-energy eigenspace for $\tilde{H}(x)$. This yields a line bundle over the parameter space which, since it is embedded in $M \times \mathbb{C}$, has a natural Hermitian connection, studied, e.g., by Bott and Chern.⁷ In this connection, one transports a vector β_0 along a curve $C(t)$ so that $\beta(t)$ obeys $\langle \beta(t + \delta t), \beta(t) \rangle = 1 + O((\delta t)^2)$. I claim that $\eta(s)$ precisely obeys this condition; formally, one can argue that

$$\begin{aligned} \left(\eta(s), \frac{d\eta(s)}{ds} \right) &= \lim_{T \rightarrow \infty} \left(\eta(s), \frac{d\eta_T(s)}{ds} \right) \\ &= \lim_{T \rightarrow \infty} (\eta(s), -iTH(s)\eta_T(s)) = 0 \end{aligned}$$

by (4). One can give a rigorous proof just using the convergence of η without worrying about the question of convergence of derivatives.⁸ Thus the connection given by the adiabatic theorem [when $E(s) \equiv 0$] is precisely the conventional one for embedded Hermitian bundles and γ is the conventional integral of the curvature which is just the Chern class of the connection. In particular, γ only depends on the X_x 's, not other aspects of $\tilde{H}(x)$. This means that one has a simple compact formula⁹ for γ :

$$\gamma(C) = \int_S V,$$

where S is any oriented surface in M with $\partial S = C$ and V can be defined in terms of an arbitrary smooth choice,¹⁰ $\varphi(x)$, of unit vectors in X_x by

$$V = i(d\varphi, d\varphi), \quad (5)$$

which is shorthand for the two-form¹¹

$$V = \sum_{i < j} \text{Im}(\partial\varphi/\partial x_i, \partial\varphi/\partial x_j) dx_i \wedge dx_j,$$

written in terms of local coordinates. The formula that Berry used has the advantage over (5) of being manifestly invariant under phase changes of $\varphi(x)$,¹² but it appears to depend on details of $\tilde{H}(x)$ and not just on the spaces X_x . Moreover, since it has a sum over intermediate states, it could be difficult to compute in general, although in his examples it is easy to compute, since the sum over intermediate state in these examples is finite. Even in these examples, (5) can be very easy to use in computations.¹³

Equation (5) shows that $V = 0$ if one can choose the $\varphi(x)$ to be all simultaneously real. Thus, the phenomena we discuss here are only present

in magnetic fields or some other condition producing a nonreal Hamiltonian.

As an example of significance below which is also considered by Berry,¹ let $M = \mathbb{R}^3 \setminus \{0\}$ and given $x \in M$, let $H(x) = \hat{x} \cdot \vec{L}$ where L is a spin- S spin on C^{2S+1} . Then all eigenvalues are nondegenerate and we can, for each $m = -S, -S+1, \dots, S$, compute a $V_m(\hat{x})$ associated to the eigenvalue $|\hat{x}|m$. By rotation covariance, $V_m(\hat{x})$ must be a function of $|\hat{x}|$ times the area form $A(\hat{x})$ on the sphere of radius \hat{x} . Thus, we need only compute V_m at $\hat{x} = (0, 0, z)$. If $|m\rangle$ is the vector with $L_z|m\rangle = m|m\rangle$, then for \hat{x} near $(0, 0, z)$, we can take

$$\varphi(\hat{x}) = \exp\left[i\left(\frac{x}{z} L_y - \frac{y}{z} L_x\right) + O(x^2 + y^2)\right] |m\rangle$$

so that

$$\begin{aligned} d\varphi &= iz^{-1}[dx L_y - dy L_x] |m\rangle, \\ (d\varphi, d\varphi) &= z^{-2} dx \wedge dy \langle m | [L_y, L_x] | m \rangle \\ &= iz^{-2} m dx \wedge dy, \end{aligned}$$

and thus, returning to general x ,

$$V_m(\hat{x}) = m |\hat{x}|^{-2} A(\hat{x}). \quad (6)$$

In particular, if S is any sphere¹⁴ about the origin,

$$(2\pi)^{-1} \int_S V_m(\hat{x}) = 2m \quad (7)$$

is an integer. This is no coincidence: If C is a clockwise circuit around the equator of the sphere S , which breaks up S into two hemispheres S_+, S_- with $\partial S_\pm = \pm C$, then

$$\exp[i\gamma(C)] = \exp(i \int_{S_+} V) = \exp(-i \int_{S_-} V)$$

so that $\int_S V$ must be 2π times an integer. We therefore see the familiar quantization of the integral of the Chern class, V , of the bundle, as a consistency condition on the holonomy, a standard fact.

Thouless *et al.*, in their deep analysis of the quantized Hall effect,⁶ considered a band of a two-dimensional solid in magnetic field, so that for each k in T^2 , the Brillouin zone, the corresponding band energy is nondegenerate. If $\varphi(k)$ is the corresponding eigenvector, then $(i/2\pi) \times \int_{T^2} (d\varphi, d\varphi)$ is an integer, the TKN^2 integer of the band. Using the source¹⁵ analogy of Berry,¹ we can "interpret" these integers. Suppose the band under consideration is the n th; and suppose an arbitrary smooth interpolation, $\tilde{H}(k)$,¹⁶ of the band Hamiltonian $H(k)$ is given from the surface of the torus into the solid torus \tilde{T} , i.e., $\tilde{H}(k)$ is defined for \tilde{k} in \tilde{T} and equals $H(k)$ on the surface.

The Wigner-von Neumann theorem¹⁷ says that generically,¹⁸ the n th band is only degenerate for isolated points $\{p_i\}_{i=1}^l$ in \tilde{T} . One can define V on \tilde{T} with these points removed, and since $dV=0$, Gauss' theorem assures us that the integral of V over the torus is just the same as its integral over little spheres about the degeneracies. Each sphere has a "charge" associated with it which is $\frac{1}{2}q_i$, with q_i an integer, and $\sum q_i$ is the TKN² integer.

It is worthwhile to expand slightly on this picture.¹⁹ Consider a matrix family $M(\tilde{x})$, depending smoothly on these parameters. If all eigenvalues of $M(\tilde{x}_0)$ are nondegenerate, we say that \tilde{x}_0 is a *regular point*. \tilde{x}_0 is a *normal singular point* if and only if (i) only one eigenvalue is degenerate and its multiplicity is 2; (ii) the degeneracy is removed to first order for any line through \tilde{x}_0 . If 0 is a normal singular point and P is the projection onto the degenerate eigenvalues of $M(0)$, then for \tilde{x} near zero,

$$PM(\tilde{x})P - PM(0)P = \tilde{a} \cdot \tilde{x} P + \tilde{B} \cdot \tilde{x} + O(x^2)$$

where \tilde{B} is a vector of traceless operators on P . Picking a basis for the range of P ,¹⁹ we can write $\tilde{B} \cdot \tilde{x} = \tilde{\sigma} \cdot C \tilde{x}$ where C is a 3×3 matrix and $\tilde{\sigma}$ are the usual Pauli matrices in the basis. The condition on removal of degeneracy says that $\det(C) \neq 0$. The Hermiticity of M implies $\det(C)$ is real. I call the sign of $\det(C)$ the *signature* $\sigma(\tilde{x}_0)$ of the normal singular point. With use of a deformation argument and the example discussed above (with spin $S=\frac{1}{2}$), it is not hard to show that if the n th and $(n+1)$ st levels are degenerate at x_0 , and if V_j is the Chern class associated to the j th level, then for a small sphere S about x_0 , we have that

$$(2\pi)^{-1} \int_S V_{n+1} = \sigma(x_0); \quad (2\pi)^{-1} \int_S V_n = -\sigma(x_0).$$

If M is a smooth, regular matrix family on T^2 and \tilde{M} is a smooth interpolation to \tilde{T} with only normal singular points,¹⁸ then the n th TKN² integer is exactly equal to a weighted sum of singular points: Points where the n th level is nondegenerate get weight zero, those where it is degenerate with the next lower level get weight $\sigma(p)$ where σ is the signature of p , and those where it is degenerate with the next higher level get weight $-\sigma(p)$.²⁰

I conclude with a mathematical remark: I have shown how the Chern integers associated to certain line bundles can be understood in terms of singularities of interpolations of the bundle. It would be interesting to extend this picture to gen-

eral vector bundles.

It is a pleasure to thank D. Robinson and N. Trudinger for the hospitality of the Australian National University, where this work was done, M. Berry for telling me of his work, and B. Souillard for the remark (via M. Berry) that there must be a connection between Berry's work and that of TKN². This research was partially supported by the National Science Foundation through Grant No. MCS-81-20833.

¹M. V. Berry, to be published; see also Proceedings of the Como Conference on Chaos, 1983 (to be published).

²See T. Kato, J. Phys. Soc. Jpn. **5**, 435-439 (1950); A. Messiah, *Quantum Mechanics* (North-Holland, Amsterdam, 1962), Vol. 2.

³There are, in the infinite-dimensional case, also domain conditions if $H(t)$ is unbounded. We ignore these in our discussion. For purposes of this note, one can think of finite-dimensional cases.

⁴Since $\int f E(s/T) ds - T \int f E(s) ds$ will be very large in the limit $T \rightarrow \infty$ unless $\int f E(s) ds = 0$, it may be difficult to set up the experiment in such a way that the "dynamic phase" $\int f E(s/T) ds$ does not wash out $\gamma(C)$.

⁵See T. Eguchi *et al.*, Phys. Rep. **66**, 213 (1980); Y. Choquet-Bruhat *et al.*, *Analysis Manifolds and Physics* (North-Holland, Amsterdam, 1983).

⁶D. Thouless, M. Kohmoto, M. Nightingale, and M. den Nijs, Phys. Rev. Lett. **49**, 405 (1982), designated as TKN².

⁷R. Bott and S. Chern, Acta Math. **114**, 71 (1965).

⁸Pick f a smooth function of compact support and compute

$$\int f(s) \left(\eta(s), \frac{d\eta}{ds} \right) ds - T \int f(s) \left(\eta_T(s), \frac{d\eta}{ds} \right) ds.$$

Integrate by parts and get one term $(d\eta_T/ds, \eta)$ which is zero by (3) and (4) and one term

$$T \int f'(s) (\eta_T, \eta) - \int f'(s) (\eta, \eta) ds - \int f'(s) ds = 0.$$

⁹While one gets (5) by appealing to the abstract theory, it is quite easy to compute. If $\eta(t) = e^{i\theta(t)} \varphi(t)$, then $(\eta(s), d\eta/ds) = 0$ becomes $(\varphi(t), d\varphi/dt) + i\dot{\theta} = 0$, so that $\gamma(C) = \int_C i(d\varphi, d\varphi) - \int_S i(d\varphi, d\varphi)$ by Stokes' theorem.

¹⁰If M has "holes," it may not be possible to choose φ globally, but since (5) is phase invariant, one need only make a choice in a neighborhood of any given point. If S is the image of a disk (as it typically is), one can always make a global choice.

¹¹This formula appears as Eq. (7b) in Ref. 1 but only in passing; surprisingly, it is not used again. For example, with it, the one-page calculation in Ref. 1 that $dV=0$ is trivial from $d^2=0$.

¹²It is an easy calculation that (5) is invariant under $\varphi(r) \rightarrow e^{i\alpha(r)} \varphi(r)$ since the normalization condition

$(\varphi, \varphi) = 1$ implies that $(\varphi, d\varphi) + (d\varphi, \varphi) = 0$. Avron, Seiler, and Simon [J. Avron, R. Seiler, and B. Simon, to be published; see also J. Avron *et al.*, Phys. Rev. Lett. **51**, 51–53 (1983)] give a simple manifestly phase-invariant form for V ; viz. $V = \frac{1}{2}i\text{Tr}(dP P dP)$ where $P(x)$ is the orthogonal projection onto X_x .

¹³I emphasize that (5) holds for the factor $\gamma(C)$ in (2) even if $E(s) \neq 0$.

¹⁴Since $dV = 0$ (i.e., $\text{div } \tilde{V} = 0$ if $V = \tilde{V}_x dx \wedge dy + \tilde{V}_y dz \wedge dx + \tilde{V}_z dy \wedge dz$), away from $\mathbf{x} = 0$, (7) holds for any surface S surrounding 0.

¹⁵Since Berry is talking about integrating (e, de) along curves which he makes analogous to a vector potential, he talks about magnetic monopoles. Since we only care about (de, de) whose dual is divergenceless away from degeneracies, we do not use the magnetic monopole language. Since the dual may not have zero curl, electrostatic language is not appropriate. Since the sources have a sign, we still use the phrase “charge” for the coefficient of the delta function in $d(de, de)$ at singularities.

¹⁶One can show (see Avron, Seiler, and Simon, and Avron *et al.*, Ref. 12) that such interpolations exist; indeed, that in the finite-dimensional case, the set of such interpolations is a dense open set in the set of all interpolations.

¹⁷J. von Neumann and E. Wigner, Phys. Z. **30**, 467 (1929); see also J. Avron and B. Simon, Ann. Phys. (N.Y.) **110**, 85–110 (1978).

¹⁸These things will be discussed further in Avron, Seiler, and Simon, Ref. 12.

¹⁹Changing basis multiplies C by a unitary and so $\sigma(x_0)$ is independent of basis. It does depend on an orientation of R^3 (order of x, y, z) but so do the TKN² integers. Everything (σ , the TKN² integers, the sign of spanning surfaces for C) changes sign under change of orientation.

²⁰The fact that the sum of the TKN² integers is zero in the finite-dimensional case is made particularly transparent by these relations. The number of points of degeneracies of level n and level $n+1$ is at least the absolute value of the sum of the first n TKN² integers.

Appearance of Gauge Structure in Simple Dynamical Systems

Frank Wilczek and A. Zee^(a)

Institute for Theoretical Physics, University of California, Santa Barbara, California 93106

(Received 9 April 1984)

Generalizing a construction of Berry and Simon, we show that non-Abelian gauge fields arise in the adiabatic development of simple quantum mechanical systems. Characteristics of the gauge fields are related to energy splittings, which may be observable in real systems. Similar phenomena are found for suitable classical systems.

PACS numbers: 03.65.Bz, 11.15.-q

Gauge fields, both Abelian and non-Abelian, figure prominently in modern theories of fundamental interactions. They also arise naturally in many geometrical contexts, and are central to much of modern mathematics. In this note we point out that gauge fields appear in a very natural way in ordinary quantum mechanical problems, whose initial formulation has no apparent relationship to gauge fields. We discuss some simple model problems in detail, and sketch in a general way how observable consequences of the gauge structures might be extracted for real physical systems. Finally, analogous behavior for classical oscillators is described.

It is, of course, potentially significant for models of elementary particles that gauge fields can arise "from nowhere," but we shall not attempt specific speculations along that line here.

Adiabatic problem.—Consider problems of the following general type: We are given a family of Hamiltonians $H(\vec{\lambda})$ depending continuously on parameters $\vec{\lambda}$, all of which have a set of n degenerate levels. By a simple renormalization of the energies, we can suppose that these levels are at $E = 0$. Such degeneracies typically will occur when for each fixed value of the $\vec{\lambda}$ there is a symmetry; however, there need not be a single symmetry which is valid for all $\vec{\lambda}$. For example, the symmetry might be rotation around an axis whose direction is specified by $\vec{\lambda}$. More generally, the symmetry group H responsible for the degeneracy is embedded in a larger group G in a $\vec{\lambda}$ -dependent way.

By the reasoning leading to the usual adiabatic theorem,¹ if the parameters are slowly varied from an initial value $\vec{\lambda}_i$ to some final value $\vec{\lambda}_f$ over a long time interval T , and the given space of degenerate levels does not cross other levels, then solutions of

$$H(\vec{\lambda}_i)\psi = 0 \quad (1)$$

are mapped onto solutions of

$$H(\vec{\lambda}_f)\psi = 0 \quad (2)$$

by solving the time-dependent Schrödinger equation

$$i\partial\psi/\partial t = H(\vec{\lambda}(t))\psi \quad (3)$$

with the boundary conditions $\vec{\lambda}(0) = \vec{\lambda}_i$, $\vec{\lambda}(T) = \vec{\lambda}_f$.

If $\vec{\lambda}_i = \vec{\lambda}_f$, so that the initial and final Hamiltonians are identical, then it becomes possible to formulate a more refined question: Given that the n degenerate levels are mapped back onto themselves by adiabatic development, is this mapping a non-trivial transformation? We find that it is, and that to describe such transformations gauge fields are the appropriate tool.

For $n = 1$, a single level, the mapping is a simple phase multiplication, or for real wave functions, a sign. These situations, corresponding to $U(1)$ or Z_2 gauge fields, were discussed by Berry² and by Simon.³

In the problem above, choose an arbitrary smooth set of bases $\psi_a(t)$ for the various spaces of degenerate levels, so that

$$H(\vec{\lambda}(t))\psi_a(t) = 0. \quad (4)$$

Such a smooth choice can always be made locally, which is sufficient for our purposes. Let us write for the solutions of the Schrödinger problem (3), with the initial condition $\eta_a(0) = \psi_a(0)$,

$$\eta_a(t) = U_{ab}(t)\psi_b(t). \quad (5)$$

In writing (5) we have assumed the adiabatic limit, which can be justified to a sufficient degree of accuracy. Our task is to determine $U(t)$. We demand that the $\eta_a(t)$ remain normalized, so that

$$0 = (\eta_b, \dot{\eta}_a) = (\eta_b, \dot{U}_{ac}\psi_c) + (\eta_b, U_{ac}\dot{\psi}_c) \quad (6)$$

which leads, in an evident notation, to the equation

$$(U^{-1}\dot{U})_{ba} = (\psi_b, \dot{\psi}_a) \equiv A_{ab}. \quad (7)$$

We will show that A , an anti-Hermitian matrix, plays the role of a gauge potential. Equation (7) is solved in terms of path-ordered integrals by

$$U(t) = P \exp \int_0^t A(\tau) d\tau. \quad (8)$$

It is remarkable that A depends only on the geometry of the space of degenerate levels. The specific form of A_{ab} computed from (7) depends, of course, upon the choice of bases $\psi_a(t)$. If one makes a different choice

$$\psi'(t) = \Omega(t)\psi(t), \quad (9)$$

then the A fields transform as

$$A'(t) = \dot{\Omega}\Omega^{-1} + \Omega A \Omega^{-1}, \quad (10)$$

i.e., as proper gauge potentials. As for ordinary gauge potentials, the path-ordered integral (8) around a closed loop transforms in a simple way under the gauge transformation (9), like (10) but without the inhomogeneous term. In particular, its eigenvalues are gauge invariant.

More generally, we can define the gauge potential A_μ everywhere on M , the space coordinatized by

$$R = \exp(i\theta_n T_{n,n+1}) \cdots \exp(i\theta_2 T_{2,n+1}) \exp(i\theta_1 T_{1,n+1}).$$

The embedding of the relevant symmetry group $\text{SO}(n)$ in $\text{SO}(n+1)$ varies with time. The parameter space M is, of course, the coset space $\text{SO}(n+1)/\text{SO}(n) = S^n$. A simple evaluation of Eq. (11) gives the non-Abelian gauge potential and field strength

$$A_\mu = \pi R^{-1}(\partial R/\partial\theta^\mu)\pi, \quad (13a)$$

$$-F_{\mu\nu} = \pi R^{-1} \frac{\partial R}{\partial\theta^\mu} (1 - \pi) R^{-1} \frac{\partial R}{\partial\theta^\nu} \pi - (\mu \leftrightarrow \nu), \quad (13b)$$

where π represents projection onto the first n components. Note that left and right projection π of the pure gauge $R^{-1}\partial_\mu R$ gives us a nontrivial $\text{SO}(n)$ gauge field. For $n=3$ we find explicitly

$$\begin{aligned} A_1 &= 0, & A_2 &= \sin\theta_1 T_{12}, \\ A_3 &= \sin\theta_1 \cos\theta_2 T_{13} + \sin\theta_2 T_{23} \end{aligned} \quad (14)$$

leading to the field strengths

$$\begin{aligned} F_{12} &= \cos\theta_1 T_{12}; & F_{12}^0 &= T_{12}, \\ F_{23} &= \cos^2\theta_1 \cos\theta_2 T_{23}; & F_{23}^0 &= T_{23}, \\ F_{13} &= \cos\theta_1 \cos\theta_2 T_{13}; & F_{13}^0 &= T_{13}. \end{aligned} \quad (15)$$

Since the metric structure

$$ds^2 = d\theta_1^2 + \cos^2\theta_1 d\theta_2^2 + \cos^2\theta_1 \cos^2\theta_2 d\theta_3^2$$

is diagonal we can define

$$F_{ij}^0 = -(g^{ij}g^{kl})^{1/2} F_{kl}, \quad (16)$$

the Cartesian tensor. As might have been anticipated from the simplicity of the starting Hamiltonians (13), the gauge structure is quite simple. In fact, the rotations induced by the ordered integrals (8) amount to parallel transport of tangent vectors to

the parameters $\vec{\lambda} = \{\lambda^1, \dots, \lambda^\mu, \dots\}$. Explicitly,

$$A_\mu^T = (\psi, \partial\psi/\partial\lambda^\mu). \quad (11)$$

The ordered integral

$$U(t) = P \exp \int_0^t A_\mu(\vec{\lambda}(t)) d\lambda^\mu \quad (12a)$$

depends only on the path and not on its parametrization. In particular, for a closed path on M one obtains the Wilson loop

$$U = P \exp \oint A_\mu d\lambda^\mu. \quad (12b)$$

As a simple illustration of the preceding framework, consider the generic example of a system with $(n+1)$ levels, of which n levels are degenerate (at zero energy by normalization). Let the Hamiltonian be $H = R(t)H_0R^{-1}(t)$. Here H_0 denotes an $(n+1)$ -dimensional matrix with the entries $(H_0)_{ij} = 0$ unless $i=j=n+1$ and $R(t) = R(\theta(t))$ is the rotation

the sphere, with the obvious identifications. Nevertheless, this very fact shows that the example involves truly non-Abelian gauge structure.

The example can obviously be generalized. With the Hamiltonian suitably parametrized on the homogeneous space G/H , we can evaluate the gauge field at the "north pole," thus obtaining from Eq. (13) a simple expression in terms of the structure constant f of G :

$$F_{\mu\nu} = f_{\mu\nu a} \pi \lambda_a \pi. \quad (17)$$

(Here λ_a denote the generators of G ; those generators not in H are labeled by a Greek index.) In particular, for the potentially physical example of a Hamiltonian with three levels, two of which are degenerate (at zero energy), and parametrized on $\text{SU}(3)/\text{SU}(2) \otimes \text{U}(1) = CP_2$ we have, in the standard $\text{SU}(3)$ notation, $F_{45} = (1 + \tau_3)$, $F_{67} = (1 - \tau_3)$, $F_{47} = -F_{56} = \tau_1$, and $F_{46} = F_{57} = \tau_2$.

Stationary states.—In many real systems there are fast and slow degrees of freedom, and then one may estimate the effect of the slow variables on the fast ones in the adiabatic approximation. An important familiar example is the Born-Oppenheimer

treatment of molecules, and we shall use the terminology of this example for definiteness, although we have not investigated any possible applications in realistic detail. In this context, an important problem is to find the stationary states, which in general requires, of course, that we treat the slow variables quantum mechanically. In order that our previous discussion, where these variables were, of course, treated classically, apply fairly directly, let us first discuss this in the correspondence or quasi-classical limit.

Suppose that the nuclei can be described quasi-classically as undergoing a motion with period $2\pi/\omega$; i.e., let them be in a quantum state of the type

$$|s\rangle = \int_0^{2\pi/\omega} e^{-ip\omega\tau} |\vec{\lambda}(\tau)\rangle d\tau, \quad (18)$$

where the label $\vec{\lambda}$ is periodic with period $2\pi/\omega$ and p is an integer. States labeled by different $\vec{\lambda}$ have negligible overlap, and $e^{-iHt}|\vec{\lambda}(t_0)\rangle = |\vec{\lambda}(t_0+t)\rangle$ to an adequate approximation, where H is the Hamiltonian for the nuclear motion calculated as if the electrons followed instantaneously. Within the stated approximations $|s\rangle$ is a stationary state, $e^{-iHt}|s\rangle = e^{-ip\omega t}|s\rangle$.

Let us suppose that for any fixed $\vec{\lambda}$ there is a symmetry guaranteeing degeneracy of two electron-

ic levels, but that the symmetry cannot be defined independent of $\vec{\lambda}$, as in the previous discussion. As we have seen, the development of these levels in response to the motion of $\vec{\lambda}$ can involve nontrivial phases and mixings after a complete period of the motion of $\vec{\lambda}$. We can diagonalize the mixing matrix and thus find states which are multiplied by phases $\exp(iy_1)$, $\exp(iy_2)$, after a period. For these states, we then find the quantization condition altered to read

$$\begin{aligned} \omega' T + y_1 &= 2\pi p \\ \text{or} \\ \omega' &= p\omega - y_1\omega/2\pi. \end{aligned} \quad (19)$$

In accordance with the correspondence principle, Eq. (19) represents the energy splittings for small p . In a more general framework one would construct an effective Lagrangian for the slow variables and treat this fully quantum mechanically. The phase we have found adiabatically represents a term in the Lagrangian linear in θ , where θ is the coordinate of nuclear rotation. Such a term contributes nothing to the classical equations of motion (in line with its origin as a pure phase) but does change the quantization condition. The rotational energy is altered from $n^2/2I$, $n = \text{integer}$, to $E_n = (n - \gamma/2\pi)^2/2I$. This agrees with Eq. (19) for large quantum numbers, viz.,

$$E_{n+p_1}(\gamma_1) - E_{n+p_2}(\gamma_2) = \frac{n}{I} \left(p_1 - p_2 - \frac{\gamma_1}{2\pi} + \frac{\gamma_2}{2\pi} \right) = \omega(p_1 - p_2) - \omega \left(\frac{\gamma_1}{2\pi} - \frac{\gamma_2}{2\pi} \right). \quad (19a)$$

An example of this general framework is the phenomenon of Λ doubling.⁴

Mechanical analogs.—Simple mechanical analogs exist for many of the systems discussed above. The point is that the mechanical equation $\ddot{x} = A\dot{x}$, for A anti-Hermitian, becomes the Schrödinger equation for $\psi = \dot{x}$.

In our study of mechanical analogs, we have uncovered other phenomena which may be interpreted as noncompact gauge fields. Consider a planar harmonic oscillator in a magnetic field perpendicular to the plane:

$$\ddot{x} + B(t)\dot{x} + \mu(t)x = 0, \quad (20)$$

$$B(t) = \gamma(t) \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}. \quad (21)$$

With γ and μ slowly varying, we have the approxi-

mate solution

$$\bar{x}(t) = \begin{pmatrix} 1 \\ i \end{pmatrix} a(t) \exp\left\{ \int_{t_0}^t \omega(\tau) d\tau \right\}, \quad (22)$$

with

$$\omega^2 + \gamma\omega - \mu = 0, \quad (23)$$

$$\dot{a}/a = -\omega/(2\omega + \gamma) \quad (24)$$

(the induced electric field has been ignored).

The second equation indicates that in response to an infinitely slow cyclic variation of the parameters in the (γ, μ) plane the amplitude a gets multiplied by a nontrivial path-dependent factor. Interestingly, amplification occurs despite the arbitrary slow variation and so the relevant gauge group is $GL(1, R)$, a noncompact group. More explicitly, the factor is given by

$$\exp(\oint da/a) = \exp(\oint A) = \exp(\int F),$$

with

$$A = \left\{ \pm \frac{1}{2(\gamma^2 + 4\mu)^{1/2}} - \frac{\gamma}{2(\gamma^2 + 4\mu)} \right\} d\gamma - \frac{1}{\gamma^2 + 4\mu} d\mu \quad (25)$$

and the field strength

$$F = \mp (\gamma^2 + 4\mu)^{-3/2}. \quad (26)$$

If either γ or μ is constant the area enclosed by the closed path in the (γ, μ) plane collapses and there is no amplification. Also, note that if $\mu = \text{const}$ our system conserves $\dot{x}^2 + \mu \ddot{x}^2 = \dot{a}^2 + a^2(\omega^2 + \mu)$. In the adiabatic limit, $a^2(\omega^2 + \mu) = \text{const}$, in contrast to the standard adiabatic theorem $a^2(\omega + \frac{1}{2}\gamma) = \text{const}$ for the case $\gamma = \text{const}$.⁵

This material is based upon research supported in part by the National Science Foundation under Grant No. PHY77-27084, supplemented by funds from the National Aeronautics and Space Administration. We thank W. Kohn and R. Schrieffer for helpful comments.

¹On leave from University of Washington, Seattle, Wash. 98195. Present address: Institute for Advanced Study, Princeton, N. J. 08540.

²M. Born and V. Fock, Z. Phys. 51, 165 (1928); L. I. Schiff, *Quantum Mechanics* (McGraw-Hill, New York, 1955), p. 290.

³M. V. Berry, to be published.

⁴B. Simon, Phys. Rev. Lett. 51, 2167 (1983).

⁵A doubling is discussed somewhat in the spirit of this note by G. Wick, Phys. Rev. 73, 51 (1948).

Taking the induced electric field into account abolishes the amplification in this particular example. Nevertheless, the phenomenon discussed does arise for differential equations describing physical systems, in particular, feedback circuits with effectively imaginary resistance or variable-length pendula acted upon by Coriolis forces.

Phase Change during a Cyclic Quantum Evolution

Y. Aharonov and J. Anandan

Department of Physics and Astronomy, University of South Carolina, Columbia, South Carolina 29208

(Received 29 December 1986)

A new geometric phase factor is defined for any cyclic evolution of a quantum system. This is independent of the phase factor relating the initial- and final-state vectors and the Hamiltonian, for a given projection of the evolution on the projective space of rays of the Hilbert space. Some applications, including the Aharonov-Bohm effect, are considered. For the special case of adiabatic evolution, this phase factor is a gauge-invariant generalization of the one found by Berry.

PACS numbers: 03.65.-w

A type of evolution of a physical system which is often of interest in physics is one in which the state of the system returns to its original state after an evolution. We shall call this a cyclic evolution. An example is periodic motion, such as the precession of a particle with intrinsic spin and magnetic moment in a constant magnetic field. Another example is the adiabatic evolution of a quantum system whose Hamiltonian H returns to its original value and the state evolves as an eigenstate of the Hamiltonian and returns to its original state. A third example is the splitting and recombination of a beam so that the system may be regarded as going backwards in time along one beam and returning along the other beam to its original state at the same time.

Now, in quantum mechanics, the initial- and final-state vectors of a cyclic evolution are related by a phase factor $e^{i\phi}$, which can have observable consequences. An example, which belongs to the second category mentioned above, is the rotation of a fermion wave function by 2π rad by adiabatic rotation of a magnetic field¹ through 2π rad so that $\phi = \pm\pi$. Recently, Berry² has shown that when H , which is a function of a set of parameters R^i , undergoes adiabatic evolution along a closed curve Γ in the parameter space, then a state that remains an eigenstate of $H(R)$ corresponding to a simple eigenvalue $E_n(R)$ develops a geometrical phase γ_n which depends only on Γ . Simon³ has given an interpretation of this phase as due to holonomy in a line bundle over the parameter space. Anandan and Stodolsky⁴ have shown how the Berry phases for the various eigenspaces can be obtained from the holonomy in a vector bundle. For the adiabatic motion of spin, this is determined by a rotation angle a , due to the parallel transport of a Cartesian frame with one axis along the spin direction, which contains the above-mentioned rotation by 2π radians as a special case. The result of a recent experiment⁵ to observe Berry's phase for light can also be understood as a rotation of the plane of polarization by this angle a .

In this Letter, we consider the phase change for all cyclic evolutions which contain the three examples above as special cases. We show the existence of a phase associated with cyclic evolution, which is universal in the sense

that it is the same for the infinite number of possible motions along the curves in the Hilbert space \mathcal{H} which project to a given closed curve \hat{C} in the projective Hilbert space \mathcal{P} of rays of \mathcal{H} and the possible Hamiltonians $H(t)$ which propagate the state along these curves. This phase tends to the Berry phase in the adiabatic limit if $H(t) \equiv H[R(t)]$ is chosen accordingly. For an electrically charged system, we formulate this phase gauge invariantly and show that the Aharonov-Bohm (AB) phase⁶ due to the electromagnetic field may be regarded as a special case. This generalizes the gauge-noninvariant result of Berry that the AB phase due to a static magnetic field is a special case of his phase. This also removes the mystery of why the AB phase, even in this special case, should emerge from Berry's expression even though the former is independent of this adiabatic approximation.

Suppose that the normalized state $|\psi(t)\rangle \in \mathcal{H}$ evolves according to the Schrödinger equation

$$H(t)|\psi(t)\rangle - i\hbar(d/dt)|\psi(t)\rangle, \quad (1)$$

such that $|\psi(\tau)\rangle = e^{i\phi}|\psi(0)\rangle$, ϕ real. Let $\Pi: \mathcal{H} \rightarrow \mathcal{P}$ be the projection map defined by $\Pi(|\psi\rangle) = \{|\psi'\rangle: |\psi'\rangle = c|\psi\rangle\}$, c is a complex number. Then $|\psi(t)\rangle$ defines a curve $C: [0, \tau] \rightarrow \mathcal{H}$ with $\hat{C} \equiv \Pi(C)$ being a closed curve in \mathcal{P} . Conversely given any such curve \hat{C} , we can define a Hamiltonian function $H(t)$ so that (1) is satisfied for the corresponding normalized $|\psi(t)\rangle$. Now define $|\tilde{\psi}(t)\rangle = e^{-if(t)}|\psi(t)\rangle$ such that $f(\tau) - f(0) = \phi$. Then $|\tilde{\psi}(\tau)\rangle = |\tilde{\psi}(0)\rangle$ and from (1),

$$-\frac{df}{dt} = \frac{1}{\hbar}\langle\psi(t)|H|\psi(t)\rangle - \langle\tilde{\psi}(t)|i\frac{d}{dt}|\tilde{\psi}(t)\rangle. \quad (2)$$

Hence, if we remove the dynamical part from the phase ϕ by defining

$$\beta \equiv \phi + \hbar^{-1} \int_0^\tau \langle\psi(t)|H|\psi(t)\rangle dt, \quad (3)$$

it follows from (2) that

$$\beta = \int_0^\tau \langle\tilde{\psi}|i(d|\tilde{\psi})/dt\rangle dt. \quad (4)$$

Now, clearly, the same $|\tilde{\psi}(t)\rangle$ can be chosen for every curve C for which $\Pi(C) = \hat{C}$, by appropriate choice of

$f(t)$. Hence β , defined by (3), is independent of ϕ and H for a given closed curve \hat{C} . Indeed, for a given \hat{C} , $H(t)$ can be chosen so that the second term in (3) is zero, which may be regarded as an alternative definition of β . Also, from (4), β is independent of the parameter t of \hat{C} , and is uniquely defined up to $2\pi n$ ($n = \text{integer}$). Hence $e^{i\beta}$ is a geometric property of the unparametrized image of \hat{C} in \mathcal{P} only.

$$a_m = -a_m \langle m | \dot{m} \rangle - \sum_{n \neq m} a_n \frac{\langle m | \dot{H} | n \rangle}{E_n - E_m} \exp \left[-\frac{i}{\hbar} \int (E_m - E_n) dt \right], \quad (5)$$

where the dot denotes time derivative. Suppose that

$$\sum_{n \neq m} \left| \frac{\hbar \langle m | \dot{H} | n \rangle}{(E_n - E_m)^2} \right| \ll 1. \quad (6)$$

Then if $a_n(0) = \delta_{nm}$, the last term in (5) is negligible and the system would therefore continue as an eigenstate of $H(t)$, to a good approximation.

In this adiabatic approximation, (5) yields

$$a_m(t) = \exp \left(- \int \langle m | \dot{m} \rangle dt \right) a_m(0).$$

For a cyclic adiabatic evolution, the phase $i \int \delta \langle m | \dot{m} \rangle dt$ is independent of the chosen $|m(t)\rangle$ and Berry² regarded this as a geometrical property of the parameter space of which H is a function. But this phase is the same as (4) on our choosing $|\psi(t)\rangle = |m(t)\rangle$ in the present approximation. But β , defined by (3), does not depend on any approximation; so (4) is exactly valid. Moreover, $|\psi(t)\rangle$ need not be an eigenstate of $H(t)$, unlike in the limiting case studied by Berry. Also, the two examples below will show respectively that it is neither necessary nor sufficient to go around a (nontrivial) closed curve in parameter space in order to have a cyclic evolution, with our associated geometric phase β . For these reasons, we regard β as a geometric phase associated with a closed curve in the projective Hilbert space and not the parameter space, even in the special case considered by Berry. But given a cyclic evolution, an $H(t)$ which generated this evolution can be found so that the adiabatic approximation is valid. Then β can be computed with the use of the expression given by Berry in terms of the eigenstates of this Hamiltonian.

We now consider two examples in which the phase β emerges naturally and is observable, in principle, even though the adiabatic approximation is not valid. Suppose that a spin- $\frac{1}{2}$ particle with a magnetic moment is in a homogeneous magnetic field \mathbf{B} along the z axis. Then the Hamiltonian in the rest frame is $H_1 = -\mu B \sigma_z$, where

$$\sigma_z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}.$$

Also,

$$|\psi(0)\rangle = \begin{pmatrix} \cos(\theta/2) \\ \sin(\theta/2) \end{pmatrix}$$

Consider now a slowly varying $H(t)$, with $H(t)|n(t)\rangle = E_n(t)|n(t)\rangle$, for a complete set $\{|n(t)\rangle\}$. If we write

$$|\psi(t)\rangle = \sum_n a_n(t) \exp \left[-\frac{i}{\hbar} \int E_n dt \right] |n(t)\rangle,$$

and use (1) and the time derivative of the eigenvector equation,⁷ we have

so that

$$|\psi(t)\rangle = \exp(i\mu B t \sigma_z / \hbar) |\psi(0)\rangle = \begin{pmatrix} \exp(i\mu B t / \hbar) \cos(\theta/2) \\ \exp(-i\mu B t / \hbar) \sin(\theta/2) \end{pmatrix},$$

which corresponds to the spin direction being always at an angle θ to the z axis. This evolution is periodic with period $\tau = \pi \hbar / \mu B$. Then from (3), for each cycle, $\beta = \pi(1 - \cos \theta)$, up to the ambiguity of adding $2\pi n$. Hence, β is $\frac{1}{2}$ of the solid angle subtended by a curve traced on a sphere, by the direction of the spin state, at the center. This is like the Berry phase except that in the latter case (1) the solid angle is subtended by a curve traced by the magnetic field $\mathbf{B}'(t)$ which is large [i.e., $\mu B' / \hbar \gg \omega$, the frequency of the orbit of $\mathbf{B}'(t)$] so that the adiabatic approximation is valid, and (2) $|\psi(t)\rangle$ is assumed to be an eigenstate of this Hamiltonian. Indeed, we may substitute such a Hamiltonian for the above H_1 or add it to H_1 with $\omega = 2\mu B / \hbar$, without changing β , in this approximation. The spin state will also move through the same closed curve in the projective Hilbert space as above if the magnetic field $\mathbf{B} = (B_0 \cos \omega t, -B_0 \sin \omega t, B_3)$ with $\cot \theta = (B_3 - \hbar \omega / 2\mu) / B_0$, where $B_0 \neq 0$.⁸ And β is the same for all such Hamiltonians. This illustrates the statement earlier that β is the same for all curves C in H with the same $\hat{C} \equiv \Pi(C)$. Also, β may be interpreted as arising from the holonomy transformation, around the closed curve on the above sphere traced by the direction of the spin state, due to the curvature on this sphere,⁴ which is a rotation. By varying appropriately a magnetic field applied to the two arms of a neutron interferometer with polarized neutrons, it is possible to make the dynamical part of β [the last term in (3)] the same for the two beams.^{2,4} Then the phase difference between the two beams is just the geometrical phase, which is observable in principle, from the interference pattern, even when the magnetic field is varied nonadiabatically. In particular, a phase difference of $\pm \pi$ rad would correspond to a 2π -rad rotation of the fermion wave function, which is thus observable.

As our second example, suppose that the magnetic field is $\mathbf{B}(t) = \mathbf{B}_0 + \mathbf{B}_1(t)$, where \mathbf{B}_0 is constant and $\mathbf{B}_1(t)$ rotates slowly in a plane containing \mathbf{B}_0 with $|\mathbf{B}_1(t)|$

$= |\mathbf{B}_0|$. Suppose that at time t the angle between \mathbf{B}_1 and \mathbf{B}_0 is $\pi - \theta(t)$ and the spin state $|\psi(t)\rangle$ is in an approximate eigenstate of $H(t) = \mu\mathbf{B} \cdot \boldsymbol{\sigma}$, where $\boldsymbol{\sigma}$ are the Pauli spin matrices. For $0 \leq \theta \ll 1$, the adiabatic condition (6) gives $0 \leq -\hbar\dot{\theta}/\mu B_0 \theta \ll 1$, assuming $\dot{\theta} \leq 0$. Hence $\theta \gg \theta_0 \exp(-\mu B_0 t/\hbar) > 0$. So θ can never become zero. That is, if $\mathbf{B}(T) = \mathbf{0}$ for some T then the adiabatic approximation, as defined above, cannot be satisfied, regardless of how slowly $\mathbf{B}_1(t)$ rotates. However, because of conservation of angular momentum, $|\psi(t)\rangle$ remains an eigenstate of $H(t)$ even at $t = T$. But if θ changes monotonically then a level crossing occurs at the point of degeneracy ($\mathbf{B} = \mathbf{0}$) so that the energy eigenvalue corresponding to $|\psi(t)\rangle$ changes sign at $t = T$. For each rotation of \mathbf{B}_1 by 2π rad, $|\psi\rangle$ rotates by π rad, so that the system returns to its original state after two rotations of $\mathbf{B}(t)$. For this cyclic evolution, our $\beta = \pi$ which can be seen from the fact that a spin- $\frac{1}{2}$ particle acquires a phase π during a rotation, or that the curve \hat{C} on the projective Hilbert space, which is a sphere, is a great circle, subtending a solid angle 2π at the center.

$$\frac{df}{dt}(t) = \langle \tilde{\psi}(t) | \frac{d}{dt} - \frac{q}{\hbar} \hat{A}_0(t) | \tilde{\psi}(t) \rangle - \frac{1}{\hbar} \langle \psi(t) | H_k(t) | \psi(t) \rangle. \quad (7)$$

We consider now a cyclic evolution so that

$$|\psi(\tau)\rangle = e^{i\phi} \exp \left(-\frac{iq}{\hbar} \int_0^\tau \hat{A}_0 dt \right) |\psi(0)\rangle.$$

Choose $f(t)$ so that $\phi = f(\tau) - f(0)$. Then

$$|\tilde{\psi}(\tau)\rangle = \exp \left(-i \frac{q}{\hbar} \int_0^\tau \hat{A}_0 dt \right) |\tilde{\psi}(0)\rangle.$$

So we now define the gauge-invariant generalization of (3) as

$$\beta \equiv \phi + \frac{1}{\hbar} \int_0^\tau \langle \psi(t) | H_k(t) | \psi(t) \rangle dt, \quad (8)$$

which on use of (7) gives

$$\beta = \int_0^\tau \langle \tilde{\psi}(t) | i \frac{d}{dt} - \frac{q}{\hbar} \hat{A}_0(t) | \tilde{\psi}(t) \rangle dt. \quad (9)$$

Here, $|\tilde{\psi}(\tau)\rangle$ is obtained by parallel transport of $|\tilde{\psi}(0)\rangle$, with respect to the electromagnetic connection, along the congruence of lines parallel to the time axis. We could have chosen, instead, any other congruence of paths from $t=0$ to $t=\tau$ in our definition of ϕ and therefore $|\tilde{\psi}(\tau)\rangle$. This would correspondingly change β , which therefore depends on the chosen congruence. But, again, β is independent of ϕ and $H(t)$ for all the motions in \mathcal{H} that project to the same closed curve \hat{C} in \mathcal{P} , for a given

This example is similar to Berry's phase in that $|\psi(t)\rangle$ is always an eigenstate of $H(t)$, even though Berry's prescription cannot be applied here because of the crossing of the point of degeneracy at which the adiabatic approximation breaks down.

Consider now a system with electric charge q for which $H = H_k(\mathbf{p} - (q/c)\hat{\mathbf{A}}(t), R_i) + q\hat{A}_0(t)$ in (1). Here, $\langle \mathbf{x} | \hat{A}_\mu(t) | \psi(t') \rangle = A_\mu(\mathbf{x}, t)\psi(\mathbf{x}, t')$, where $A_\mu(\mathbf{x}, t)$ is the usual electromagnetic four-potential, and R_i are some parameters. Under a gauge transformation,

$$|\psi(t)\rangle \rightarrow \exp[i(q/c)\hat{A}(t)] |\psi(t)\rangle,$$

$$\hat{A}_0(t) \rightarrow \hat{A}_0(t) - c^{-1} \partial \hat{A}(t) / \partial t,$$

and

$$H_k(t) \rightarrow \exp[i(q/c)\hat{A}(t)] H_k(t) \exp[-i(q/c)\hat{A}(t)].$$

As before, define $|\tilde{\psi}(t)\rangle = e^{-if(t)} |\psi(t)\rangle$. If we require that $|\tilde{\psi}\rangle$ undergo the same gauge transformation as $|\psi(t)\rangle$, $f(t)$ is gauge invariant. Then, from (1),

chosen congruence. Both β and ϕ , which satisfies

$$e^{-i\phi} = \langle \psi(t) | \exp \left(-\frac{iq}{c} \int_0^\tau \hat{A}_0 dt \right) | \psi(0) \rangle,$$

are gauge invariant. In the adiabatic limit, $|\tilde{\psi}(t)\rangle$ can be chosen to be an eigenstate of $H_k(t)$ and (9) is then a gauge-invariant generalization of the Berry phase.

We illustrate this by means of the AB effect.⁶ Berry has obtained the AB phase from the gauge-noninvariant expression (4) with $|\tilde{\psi}(t)\rangle$ an eigenstate of $H(t)$, for a stationary magnetic field, in a special gauge.⁹ But a gauge can be chosen so that the AB phase is included in the dynamical phase instead of the geometrical phase (4). Also, in general, there is no cyclic evolution in an AB experiment. But our β defined by Eq. (8) or (9) is gauge invariant and includes the AB phase in the special case to be described now.

Suppose that a charged-particle beam is split into two beams at $t=0$ which, after traveling in field-free regions, are recombined so that they have the same state at $t=t$. It is assumed here that the splitting and the subsequent evolution of the two beams occur under the action of two separate Hamiltonians. This is possible if we restrict ourselves to the Hilbert space of a subset of the degrees of freedom of a given system, as in the example considered by Aharonov and Vardi.¹⁰ This belongs to the third example of a cyclic evolution mentioned at the beginning of this Letter. The wave function of each beam

at $t = \tau$, assuming that it has a fairly well defined momentum, is

$$\psi(x, \tau) = \exp\left(-\frac{i}{\hbar} \int_0^\tau E_i dt\right) \exp\left(-\frac{iq}{c} \int_{\gamma} A_\mu dx^\mu\right) \exp\left(\frac{i}{\hbar} \int_{\gamma} p \cdot dx\right) \psi(x, 0), \quad i = 1 \text{ or } 2,$$

where γ is a space-time curve through the beam and p represents the approximate kinetic momentum of the beam. Hence on using (8), we have

$$\beta = -\frac{q}{c} \oint_{\gamma} A_\mu dx^\mu + \frac{1}{\hbar} \oint_{\gamma} p \cdot dx, \quad (10)$$

where γ is the closed curve formed from γ_1 and γ_2 . But this is only an approximate treatment and a more careful investigation of this problem is needed.

In conclusion, we note that $\mathcal{H}^* = \mathcal{H} - \{0\}$ is a principal fiber bundle over \mathcal{P} with structure group C^* (the group of nonzero complex numbers), and the disjoint union of the rays in \mathcal{H} is the natural line bundle over \mathcal{P} whose fiber above any $p \in \mathcal{P}$ is p itself. Then, clearly, β , given by (4), arises from the holonomy due to a connection in either bundle such that $|\psi(t)\rangle$ is parallel transported if

$$\langle \psi(t) | (d/dt) | \psi(t) \rangle = 0, \quad (11)$$

i.e., the horizontal spaces are perpendicular to the fibers with respect to the Hilbert space inner product. Condition (11) was used by Simon³ to define a connection on a line bundle over parameter space, which is different from the above bundles. The real part of (11) says that $\langle \psi(t) | \psi(t) \rangle$ is constant during parallel transport. Since this is true also during any time evolution determined by (1), we may restrict consideration to the subbundle $\mathcal{F} = \{|\psi\rangle \in \mathcal{H}: \langle \psi | \psi \rangle = 1\}$ of \mathcal{H}^* . This \mathcal{F} is the Hopf bundle¹¹ over \mathcal{P} . Then the imaginary part of (11) defines the horizontal spaces in \mathcal{F} which determine a connection. This is the usual connection in \mathcal{F} and $e^{i\beta}$ is the holonomy transformation associated with it. If \mathcal{H} has finite dimension N then \mathcal{P} has dimension $N-1$. For $N=2$, \mathcal{P} is the complex projective space $P_1(C)$ which is a sphere with the Fubini-Study metric¹¹ on \mathcal{P} being the usual metric on the sphere. Opposite points on this sphere represent rays containing orthogonal states. Our geometric phase can then be obtained from the holono-

my angle α associated with parallel transport around a closed curve on this sphere like in Ref. 4.

It is a pleasure to thank Don Page for suggesting the relevance of the Hopf bundle and the Fubini-Study metric to this work.

¹Y. Aharonov and L. Susskind, Phys. Rev. **158**, 1237 (1967).

²M. V. Berry, Proc. Roy. Soc. London, Ser. A **392**, 45 (1984).

³B. Simon, Phys. Rev. Lett. **51**, 2167 (1983).

⁴J. Anandan and L. Stodolsky, Phys. Rev. D **35**, 2597 (1987).

⁵R. Y. Chiao and Y.-S. Wu, Phys. Rev. Lett. **57**, 933 (1986); A. Tomita and R. Y. Chiao, Phys. Rev. Lett. **57**, 937 (1986).

⁶Y. Aharonov and D. Bohm, Phys. Rev. **115**, 485 (1959).

⁷See, for example, L. I. Schiff, *Quantum Mechanics* (McGraw-Hill, New York, 1968), pp. 289–291.

⁸An experiment of this type has been done to measure Berry's phase ($\omega \rightarrow 0$) using nuclear magnetic resonance by D. Suter, G. Chingas, R. A. Harris, and A. Pines, to be published. One of us (J.A.) wishes to thank A. Pines for a discussion during which it was realized that the same type of experiment can be used to measure the geometric phase β introduced in the present Letter for nonadiabatic cyclic evolutions as well.

⁹In this proof, in Ref. 2, the eigenfunctions, in the absence of the electromagnetic field, are in effect assumed to be real, in order that Eq. (34) is valid. Since the coefficients of the stationary Schrödinger equation are then real, it is always possible to find real solutions. Then, for any eigenfunction belonging to a given eigenvalue to be necessarily a real function multiplied by $e^{i\lambda}$ ($\lambda = \text{const}$), it is necessary and sufficient that the eigenvalue is simple. But in our treatment of the AB effect, it is not necessary to make this assumption.

¹⁰Y. Aharonov and M. Vardi, Phys. Rev. D **20**, 3213 (1979).

¹¹See, S. Kobayashi and K. Nomizu, *Foundations of Differential Geometry* (Interscience, New York, 1969), Vol. 2.

PHYSICAL REVIEW

LETTERS

VOLUME 60

6 JUNE 1988

NUMBER 23

General Setting for Berry's Phase

Joseph Samuel and Rajendra Bhandari

Raman Research Institute, Bangalore 560080, India

(Received 30 November 1987)

It is shown that Berry's phase appears in a more general context than realized so far. The evolution of the quantum system need be neither unitary nor cyclic and may be interrupted by quantum measurements. A key ingredient in this generalization is the use of some ideas introduced by Pancharatnam in his study of the interference of polarized light, which, when carried over to quantum mechanics, allow a meaningful comparison of the phase between any two nonorthogonal vectors in Hilbert space.

PACS numbers: 03.65.Bz, 02.40.+m

Three years ago, Berry¹ made a rather perceptive and interesting observation regarding the behavior of quantum-mechanical systems in a slowly changing environment. If the system is initially in an eigenstate of the instantaneous Hamiltonian, the adiabatic theorem guarantees that it remains so. This, however, determines the state of the system only up to a phase. Berry asked the question "What is the phase of the system?" and got a somewhat unexpected answer. If the environment (more precisely, the Hamiltonian) returns to its initial state, the system also does, but it acquires an extra phase over and above the dynamical phase, which can be calculated and allowed for. This effect has been studied and measured in various contexts.

Simon² gave a simple geometrical interpretation of Berry's phase. If one regards the space of normalized states as a fiber bundle over the space of rays³ (a ray is defined as an equivalence class of states differing only in phase), then this bundle has a natural connection. This connection permits a comparison of the phases of states on two neighboring rays. Simon observed that when the dynamical phase factor is removed, the evolution of the system as determined by the Schrödinger equation is a parallel transport of the phase of the system according to this natural connection. Berry's phase is then a consequence of the curvature of this connection.

Recently, Aharonov and Anandan⁴ generalized Berry's results by giving up the assumption of adiabaticity. The key step in this work is their identification of the in-

tegral of the expectation value of the Hamiltonian as the dynamical phase. Once this dynamical phase is removed, the evolution of the phase of the system is again determined by the natural connection and one recovers Berry's phase for any cyclic evolution of the quantum system.

The purpose of this Letter is to point out that Berry's phase appears in a still more general context. The evolution of the system need be neither unitary (norm preserving) nor cyclic (returning to the original ray). This generalization is based on the work of Pancharatnam⁵ on the interference of polarized light. Carrying Pancharatnam's ideas over to quantum mechanics yields a fairly general setting for a discussion of Berry's phase. We briefly describe Pancharatnam's work before developing the subject of the present paper.

Pancharatnam posed the following question: Given two beams of polarized light, is there a natural way to compare the phases of these beams? His physically motivated answer was to cause interference of these two polarized beams and regard them as "in phase" when the resultant intensity is maximum. This provides a "connection" (a rule for the comparison of phases) between any two states of polarization which are not orthogonal. (This rule breaks down for orthogonal states. These do not interfere and the resultant intensity is insensitive to the relative phase of the two beams.) Consider three (nonorthogonal) states of polarization represented by three points 1, 2, and 3 on the Poincaré sphere. Suppose

now that 1 and 2 are "in phase" and 2 and 3 are "in phase"; then 1 and 3 are not necessarily in phase. Pancharatnam showed that the excess phase of 3 over 1 is given by half the solid angle subtended by the spherical triangle 123 at the center of the Poincaré sphere. Thus the Pancharatnam connection has curvature. While Pancharatnam's studies, both theoretical and experimental, were carried out in the 1950's, the relation between his work and Berry's was pointed out only recently by Ramaseshan and Nityananda.⁶ They, and subsequently Berry,⁷ observed that Pancharatnam's excess phase is in fact an early example of Berry's phase. A laser interferometer experiment demonstrating Pancharatnam's excess phase has earlier been reported by us.⁸ In the rest of this paper we show that carrying over Pancharatnam's ideas to quantum mechanics leads to a fruitful generalization of Berry's phase.

Consider a quantum system whose state vector $|\psi\rangle$ (an element of a Hilbert space \mathcal{H}) evolves according to the Schrödinger equation $i(d/dt)|\psi(t)\rangle = \hat{H}(t)|\psi(t)\rangle$. (\hat{H} is a linear operator, possibly non-Hermitean.) Let us define a new state vector $|\phi(t)\rangle$, which differs from $|\psi(t)\rangle$ only in that it has had a dynamical phase factor removed:

$$|\phi(t)\rangle = \exp\left[i \int_0^t h(t') dt'\right] |\psi(t)\rangle,$$

where

$$h(t') = (\psi|\psi)^{-1} \operatorname{Re}(\psi(t')|\hat{H}(t')|\psi(t')).$$

Clearly, $|\phi(t)\rangle$ satisfies the equation

$$i(d/dt)|\phi(t)\rangle = [\hat{H}(t) - h(t)]|\phi(t)\rangle.$$

Contracting this with $\langle\phi(t)|$ yields the parallel-transport law

$$\operatorname{Im}(\phi(t))| (d/dt)|\phi(t)\rangle = 0. \quad (1)$$

While this law has its origin in the Schrödinger equation, it is purely geometric, as are the considerations in the rest of this paper.

Let \mathcal{N} denote the set of normalizable states in \mathcal{H} :

$$\mathcal{N} = \{|\psi\rangle \in \mathcal{H} | (\psi|\psi) \neq 0\}.$$

Let \mathcal{R} be the space of rays: $\mathcal{R} = \mathcal{N}/\sim$, where \sim denotes that elements of \mathcal{N} which differ only by a phase are regarded as equivalent. There is a natural projection map $\pi: \mathcal{N} \rightarrow \mathcal{R}$, which maps each vector to the ray on which it lies. The triplet $(\mathcal{N}, \mathcal{R}, \pi)$ forms a principal fiber bundle over the base space \mathcal{R} [with structure group $U(1)$] and the parallel-transport law (1) defines a natural connection on this fiber bundle. A connection⁹ is an assignment of a "horizontal subspace" in the tangent space of each point in \mathcal{N} . Horizontal vectors are those that satisfy (1). Given a curve in \mathcal{R} , one can lift this curve up to \mathcal{N} so that its tangent vector is horizontal. However, the horizontal lift of a closed curve in \mathcal{R} may

be open in \mathcal{N} . This is referred to as holonomy of the connection and provides a geometric picture of Berry's phase.

Let $|\phi(s)\rangle$ be a curve in \mathcal{N} . Let $|u\rangle = (d/ds)|\phi(s)\rangle$ denote the tangent vector to this curve. Let us define¹⁰

$$A_s = \operatorname{Im}(\phi|u\rangle)/\langle\phi|\phi\rangle. \quad (2)$$

Under transformations of the kind $|\phi(s)\rangle \rightarrow \exp[i \times a(s)]|\phi(s)\rangle$ (referred to as gauge transformations), A_s transforms inhomogeneously,

$$A_s \rightarrow A_s + da/ds, \quad (3)$$

like the vector potential in electrodynamics. The parallel-transport law (1) states that A_s vanishes along the actual curve $|\phi(s)\rangle$ followed by the quantum system.

Let us first consider $|\psi(t)\rangle$, a solution of the Schrödinger equation which is cyclic, i.e., returns to the initial ray at some time τ . This defines a curve in \mathcal{N} . The "shadow" of this curve under projection map π is a closed curve in \mathcal{R} . Given the closed curve $r(s)$ in \mathcal{R} , let us ask for the curve $|\phi(s)\rangle$ in \mathcal{N} traced out by the state vector (with the dynamical phase removed). Using (1), we find that the curve is determined by the condition $A_s = 0$ along the curve. Consider the integral

$$\gamma = \oint A_s ds \quad (4)$$

along the curve $|\phi(s)\rangle$ in \mathcal{N} closed by the vertical curve joining $|\phi(\tau)\rangle$ to $|\phi(0)\rangle$. The segment $|\phi(s)\rangle$ represents the actual evolution of the system and along this, $A_s = 0$. The vertical contributes the phase difference between $|\phi(0)\rangle$ and $|\phi(\tau)\rangle$ and represents Berry's phase. However, the integral (4) is gauge invariant because of (3) and can be regarded as an integral on \mathcal{R} . With use of Stokes's theorem, γ can be expressed as

$$\gamma = \int_S F, \quad (5)$$

where S is any surface in \mathcal{R} bounded by the closed curve $r(s)$ in \mathcal{R} and F is the gauge-invariant ("field strength") two-form representing the curl of A . γ depends only on the geometric curve $r(s)$ and not on the rate at which it is transversed in time. This gives the formula for Berry's phase in a (possibly nonunitary) cyclic evolution of a quantum system.

In a general evolution, the state vector may not return to the initial ray. In order to handle this situation, we need a method of comparing states on different rays for phase. This is provided by the Pancharatnam connection. Let $|\phi_1\rangle$ and $|\phi_2\rangle$ be any two elements of \mathcal{N} which are not orthogonal. Interference of these two states by superposition yields

$$\| |\phi_1\rangle + |\phi_2\rangle \| ^2 = \langle \phi_1 | \phi_1 \rangle + \langle \phi_2 | \phi_2 \rangle + 2 \operatorname{Re}(\phi_1 | \phi_2).$$

The modulus of the resultant vector is clearly a maximum when $\langle \phi_1 | \phi_2 \rangle$ is real and positive. Under this condition, $|\phi_1\rangle$ and $|\phi_2\rangle$ are said to be "in phase." More

generally, if one writes the complex number $\langle \phi_1 | \phi_2 \rangle$ in polar form, $\rho \exp i\beta$, $\rho > 0$, then the phase difference between $|\phi_1\rangle$ and $|\phi_2\rangle$ is β . The Pancharatnam connection has a clear physical basis and is more general than the natural connection since it permits a comparison of *any* two (nonorthogonal) states for phase and not just neighboring ones. In the particular case where $|\phi_1\rangle$ and $|\phi_2\rangle$ are on neighboring rays, the Pancharatnam connection reduces to the natural connection.

We now go on to express the Pancharatnam phase difference in terms of the natural connection. In order to do this, we need to explore some geometrical properties of the ray space \mathcal{R} . We observe that \mathcal{R} has a natural metric on it, which comes from the (positive definite) inner product (\cdot) on \mathcal{H} . Since each point of \mathcal{R} is an entire equivalence class, it is convenient to take representative elements from \mathcal{N} and make sure our considerations are gauge invariant. Let $|\phi(s)\rangle$ be a curve in \mathcal{N} and $|u\rangle$ its tangent vector. Under gauge transformations, $|u\rangle$ does not transform covariantly. But its projection orthogonal to the fiber,

$$|u' \rangle = |u\rangle - |\phi\rangle [\langle \phi | u \rangle - \langle u | \phi \rangle] (2\langle \phi | \phi \rangle)^{-1}$$

is gauge covariant. $|u'\rangle$ is in fact the covariant derivative

$$|u' \rangle = (d/ds) |\phi(s)\rangle - iA_s |\phi(s)\rangle.$$

$\langle u' | u' \rangle$ is gauge invariant and can be used to define a metric on \mathcal{R} : $dl^2 = \langle u' | u' \rangle ds^2$. dl^2 is the square of the distance between points $\pi(|\phi(s)\rangle)$ and $\pi(|\phi(s+ds)\rangle)$. This metric can also be expressed with use of the density matrix $\rho = |\psi\rangle\langle\psi|$, which contains information only about the ray and not the phase. Its form is

$$dl^2 = (\text{Tr}\rho)^{-1} [\text{Tr}(d\rho d\rho) - \frac{1}{2} (\text{Tr}d\rho)^2].$$

This metric then determines geodesics in \mathcal{R} . These can be found by variation of $\int \langle u' | u' \rangle dl$, where l is an affine parameter. This yields the geodesic equation

$$\frac{D^2}{dl^2} |\phi(l)\rangle = \frac{d}{dl} |u' \rangle - iA_s |u' \rangle = 0. \quad (6)$$

Curves in \mathcal{N} which satisfy this equation project down to geodesics in \mathcal{R} . Notice that (6) is gauge covariant and so the geodesic nature is a property of the "shadow" of the curve and not the curve itself.

The importance of geodesic curves in \mathcal{R} stems from the fact that one can express the Pancharatnam phase difference as a line integral of A_s with use of the *geodesic rule*: Let $|\phi_1\rangle$ and $|\phi_2\rangle$ be any two (nonorthogonal) states in \mathcal{N} , with phase difference β according to the Pancharatnam connection. Let $|\phi(s)\rangle$ be any geodesic curve connecting $|\phi_1\rangle$ to $|\phi_2\rangle$: $|\phi(0)\rangle = |\phi_1\rangle$, $|\phi(1)\rangle = |\phi_2\rangle$. Then β is given by

$$\beta = \int A_s ds, \quad (7)$$

where A_s is given by (2).

Proof: Let $r(s)$ be a geodesic curve in \mathcal{R} joining $\pi(|\phi_1\rangle)$ to $\pi(|\phi_2\rangle)$. Consider the horizontal lift $|\tilde{\phi}(s)\rangle$ of this curve, which starts from $|\phi_1\rangle$ [$|\tilde{\phi}(0)\rangle = |\phi_1\rangle$, $\tilde{A}_s = 0$]. The geodesic equation (6) reduces to $(d^2/ds^2)|\tilde{\phi}(s)\rangle = 0$, whose solution is a straight line in \mathcal{N} . Further, $|\tilde{\phi}(s)\rangle$ is "in phase" with $|\tilde{\phi}(0)\rangle$. To see this, define $g(s) = \text{Im}(\tilde{\phi}(0)|\tilde{\phi}(s)\rangle)$. Clearly, $g(0) = 0$ and $g'(0) = 0$ since $|\tilde{\phi}(s)\rangle$ is a horizontal curve. Now $\dot{g}(s)$ can be worked out from the geodesic equation $\ddot{g}(s) = \text{Im}(\tilde{\phi}(0)|(d^2/ds^2)|\tilde{\phi}(s)\rangle) = 0$, so that $g(s)$ is identically zero along the horizontal curve $|\tilde{\phi}(s)\rangle$; hence $\langle \phi_1 | \tilde{\phi}(s)\rangle$ is real,¹¹ and so $|\tilde{\phi}(s)\rangle$ and $|\phi_1\rangle$ are "in phase." To prove (7), we simply perform a gauge transformation $|\phi(s)\rangle = \exp[i\alpha(s)]|\tilde{\phi}(s)\rangle$, where $\alpha(s)$ is chosen so that $\alpha(0) = 0$, $\alpha(1) = \beta$. Then $|\phi(s)\rangle$ is still a geodesic curve (since the geodesic equation is gauge covariant) and connects $|\phi_1\rangle$ to $|\phi_2\rangle$. Using the fact that \tilde{A}_s was zero along the horizontal curve $|\tilde{\phi}(s)\rangle$ and the behavior (3) of A_s under gauge transformations, we find that the right-hand side of (7) becomes $\int \delta(da/ds)ds = \beta$ and (7) is verified. The geodesic rule¹² (7) is the main result of this paper.

We are now ready to show how Berry's phase also appears in a noncyclic evolution. Let the state vector $|\phi(t)\rangle$ (with the dynamical phase removed as always) evolve from $|\phi(0)\rangle$, initially, to $|\phi(\tau)\rangle$. If $|\phi(\tau)\rangle$ is not orthogonal to $|\phi(0)\rangle$, it is meaningful to ask, "What is the phase difference between them?" If we use the geodesic rule, this can be expressed as in (7). Now add to this integral the quantity $\int A_s ds$ integrated along the actual curve determined by the Schrödinger equation. Because of (1), this vanishes and the phase difference γ between $|\phi(\tau)\rangle$ and $|\phi(0)\rangle$ is expressed as an integral (4) where the contour C is given by the actual evolution $|\phi(t)\rangle$ from $|\phi(0)\rangle$ to $|\phi(\tau)\rangle$ and back along any geodesic curve joining $|\phi(\tau)\rangle$ to $|\phi(0)\rangle$. This expression for γ is gauge invariant and so can be regarded as defined on the base \mathcal{R} . So γ can be expressed as the integral (5) of the two-form F over a surface bounded by the closed curve $\pi(C)$. γ is clearly a gauge-invariant quantity and measurable. It has a purely geometric origin and depends only on the geometric path that the system traces in \mathcal{R} and not on its rate of traversal.

Let us next consider a quantum system undergoing a nonunitary evolution, as happens, for example, when the system is subjected to measurements. According to the collapse postulate, the effect of the measurement on the system is described by the projection operator $P = |\psi\rangle\langle\psi|$ onto the eigenstate corresponding to the eigenvalue (the outcome of the measurement) of the operator measured. Consider a system initially in the state $|\psi_1\rangle$ on which three successive measurements are made. If the effects of these measurements are to project the state of the system onto $|\psi_2\rangle$, then onto $|\psi_3\rangle$ and back onto $|\psi_1\rangle$, the final state of the system is given

by $|\psi_1\rangle\langle\psi_1|\psi_1\rangle\langle\psi_3|\psi_2\rangle\langle\psi_2|\psi_1\rangle$. (We ignore the time evolution, i.e., set $H=0$, so that we can concentrate on the effects due to projection that we are interested in.) The final and initial states have a well-defined phase difference, given by the phase of the complex number $\langle\psi_1|\psi_3\rangle\langle\psi_3|\psi_2\rangle\langle\psi_2|\psi_1\rangle$. Using the geodesic rule, we see that the phase is given by (5), where now the surface is bounded by the geodesic triangle connecting rays 1, 2, and 3. Thus Berry's phase also appears in systems subjected to quantum measurements. For a spin- $\frac{1}{2}$ (two-state) system, this formula for γ reduces to half the solid angle subtended at the center of the Poincaré sphere by the rays 1, 2, and 3. γ can be experimentally measured.

In summary, Berry's phase appears to be more general than the context in which it was discovered by Berry, i.e., for an adiabatic, cyclic, and unitary evolution. Our discussion uses a new ingredient—the Pancharatnam connection—and contains previous work as special cases. Since Berry's phase is being studied and applied in many different contexts, this generalization may be of interest.

It is a pleasure to thank S. Ramaseshan and Rajaram Nityananda for exciting our interest in Pancharatnam's work and S. Sridhar for useful discussions.

Narasimhan and S. Ramanan, Am. J. Math. **83**, 563 (1961).

³Actually, Simon used the induced bundle on a finite-dimensional parameter space mapped injectively into the ray space. Y. Aharonov and J. Anandan, Phys. Rev. Lett. **58**, 1593 (1987), work directly with the ray space. See also Don Page, Phys. Rev. A **36**, 3479 (1987).

⁴Aharonov and Anandan, Ref. 3.

⁵S. Pancharatnam, Proc. Indian Acad. Sci. A**44**, 247 (1956), reprinted in *Collected Works of S. Pancharatnam* (Oxford Univ. Press, London, 1975).

⁶S. Ramaseshan and R. Nityananda, Curr. Sci. **55**, 1225 (1986).

⁷M. V. Berry, J. Mod. Opt. **34**, 1401 (1987). A proof of the geodesic rule is given here for two-state quantum systems.

⁸R. Bhandari and J. Samuel, Phys. Rev. Lett. **60**, 1211 (1988).

⁹S. Kobayashi and K. Nomizu, *Foundations of Differential Geometry* (Wiley, New York, 1963).

¹⁰The connection can be given by our specifying a connection one-form A on \mathcal{N} . The contraction of this one-form with a tangent vector $|u\rangle$ is denoted by A_u in the text. Note that A is a linear functional on tangent vectors over the reals (not over \mathbb{C}).

¹¹In fact, $\langle\delta(0)|\delta(s)\rangle$ is also positive if $|\delta(s)\rangle$ is the shortest geodesic connecting $|\delta(0)\rangle$ with $|\delta(1)\rangle$. Along this curve, $\langle\delta(0)|\delta(s)\rangle$ never vanishes and, since it was positive to start with, remains so.

¹²This rule breaks down for orthogonal states, since these are connected by a continuous infinity of geodesics, and this prescription does not give a definite answer.

¹M. V. Berry, Proc. Roy. Soc. London **392**, 45 (1984).

²B. Simon, Phys. Rev. Lett. **51**, 2167 (1983); M. S.

Progress of Theoretical Physics, Vol. 74, No. 3, September 1985

Effective Action for Adiabatic Process

— *Dynamical Meaning of Berry and Simon's Phase* —

Hiroshi KURATSUJI* and Shinji IIDA

Department of Physics, Kyoto University, Kyoto 606

(Received March 16, 1985)

By applying the path integral method to two interacting systems, it is shown that the specific phase Γ appearing in the quantum adiabatic process recently found by Berry and Simon is obtained as an additive action to the conventional dynamical action function. This scheme naturally gives a dynamical meaning to Berry and Simon's topological phase, which leads to a novel form of semiclassical quantization rule including the phase Γ .

Following an early implication in molecular physics,¹⁾ Berry²⁾ recently discovered a rather unexpected phenomenon on quantum adiabatic theorem: During an excursion along a closed loop C in the external parameter space the adiabatic change of a wave function gains an extra phase in addition to the conventional dynamical phase; $\psi_n(T) = \exp[i\Gamma_n(C)]\exp[-i/\hbar\int_0^T E_n(R(t))dt]\|n(R(T))\rangle$. Here Berry demonstrated that the phase $\Gamma(C)$ reflects a peculiar structure associated with the level-crossing which is inherent to the global geometric nature of the external parameter space. Subsequently Simon³⁾ gave this specific phase a topological meaning that it is nothing but the "holonomy" constructed from the vector bundle of parametrized wave functions and discussed a connection with the quantized Hall effect.⁴⁾ The topological concept such as vector bundle or connection has now become familiar in gauge field theory.⁵⁾ However it is rather surprising that the topological concept appears even in the usual non-relativistic quantum mechanics.

Although Berry and Simon's phase (B-S phase) has such an appealing feature, it is still concerned with the static aspect only, i.e., time development of external parameter space is given from the outset. Actually, the parameter space itself can be regarded as a dynamical object. For example, in the Born-Oppenheimer theory, the internuclear distance, which is frozen in the adiabatic process, should be regarded as a dynamical variable. Thus we are forced to inquire a dynamical meaning to Berry and Simon's topological phase. The purpose of this paper is to put forward an answer to this question. A similar dynamical argument was suggested by Mead and Truhlar⁶⁾ before Berry and Simon; they showed that this specific phase acquires a meaning of the effective vector potential in the Schrödinger equation for the nuclear motion in molecular collisions. However, the argument based on the Schrödinger equation is of essentially local nature and does not seem to be appropriate for describing the global character of this non-integrable phase. In order to push the global aspect forward, we adopt the path integral formulation for the bound state problem of interacting two systems. In this formulation, the B-S phase naturally arises as an additional action to the conventional action function

* Present address: Department of Mathematics and Physics, Ritsumeikan University, Kyoto 603.

induced by adiabatic process. Simultaneously this topological action is shown to modify the semiclassical quantization rule for the motion of the external system.

Effective action by path integral: Consider two interacting systems, which are described by the variables conventionally called “internal” and “collective” coordinates; q and X respectively. We adopt a Hamiltonian $\hat{H} = \hat{h}(q, X) + \hat{H}_0(P, X)$, where the internal Hamiltonian \hat{h} is assumed to depend on X and not on its conjugate momentum P . Let us consider the trace of the evolution operator $K(T) = \text{Tr}[\exp(-i\hat{H}T/\hbar)]$, which is written as

$$K(T) = \sum_n \int \langle n(X_0), X_0 | \exp[-i\hat{H}T/\hbar] | n(X_0), X_0 \rangle d\mu(X_0). \quad (1)$$

In Eq. (1) one naturally picks up the transition amplitude for the quantum process starting from the initial state of product form $|n(X_0), X_0\rangle (\equiv |n(X_0)\rangle \otimes |X_0\rangle)$ and returning via closed loops C to the same state, where $|X_0\rangle$ denotes the eigenstate of \hat{X} and $|n(X_0)\rangle$ is the eigenstate of $\hat{h}(q, X_0)$ at $X = X_0$ with eigenvalue $E_n(X_0)$. Then, with the aid of the time-discretization together with the completeness relation holding for X , we get

$$\begin{aligned} & \langle n(X_0), X_0 | \exp[-i\hat{H}T/\hbar] | n(X_0), X_0 \rangle \\ &= \int \prod_{k=1}^{N-1} d\mu(X_k) \langle n(X_0), X_0 | \exp[-i\hat{H}\varepsilon/\hbar] | X_{N-1} \rangle \cdots \langle X_1 | \exp[-i\hat{H}\varepsilon/\hbar] | n(X_0), X_0 \rangle \end{aligned} \quad (2)$$

with $\varepsilon = T/N$. Further noting the relation for $\varepsilon \approx 0$,

$$\begin{aligned} & \langle X_k | \exp[-i\hat{H}\varepsilon/\hbar] | X_{k-1} \rangle \simeq \langle X_k | \exp[-i\hat{H}_0\varepsilon/\hbar] | X_{k-1} \rangle \exp[-i\hat{h}(X_k)\varepsilon/\hbar] \\ &= \int dP_k \exp[(iP_k(X_k - X_{k-1}) - iH_0(X_k, P_k)\varepsilon)/\hbar] \exp[-i\hat{h}(X_k)\varepsilon/\hbar], \end{aligned} \quad (3)$$

Eq. (1) can be expressed as

$$K(T) = \sum_n \int T_{nn}(C) \exp\left[-\frac{i}{\hbar} S_0(C)\right] \prod_t d\mu(X_t, P_t) \quad (4)$$

with $S_0(C) (\equiv \int (P\dot{X} - H_0) dt)$ the action for the collective motion along closed loops C . $T_{nn}(C)$ is just the internal transition amplitude and given by

$$T_{nn}(C) = \langle n(X_0) | \exp[-i\hat{h}(N)\varepsilon/\hbar] \cdots \exp[-i\hat{h}(1)\varepsilon/\hbar] | n(X_0) \rangle, \quad (5)$$

i.e., the time ordered product, where $\hat{h}(k)$ denotes the internal Hamiltonian at the point $X = X_k$ on the loop C .⁶⁾ Namely, if we denote $|\phi_n(T)\rangle$ as a solution of the time-dependent Schrödinger equation; $(i\hbar\partial/\partial t - \hat{h}(q, X_t))|\phi_n(t)\rangle = 0$ with the boundary condition $|\phi_n(0)\rangle = |n(X_0)\rangle$, $T_{nn}(C)$ is written as $T_{nn}(C) = \langle n(X_0) | \phi_n(T) \rangle$.

Under the above prescription we turn to the case of the adiabatic motion where the period T is large. By inserting the completeness relation holding for the internal state on each point of external variables X_k ; $\sum_m |m_k\rangle \langle m_k| = 1$, Eq. (5) is written as

$$T_{nn}(C) = \sum_{m_1} \cdots \sum_{m_{N-1}} \langle n(X_0) | \exp[-i\hat{h}(N)\varepsilon/\hbar] | m_{N-1} \rangle$$

$$\cdots \langle m_k | \exp[-i\hat{h}(k)\varepsilon/\hbar] | m_{k-1} \rangle \cdots \langle m_1 | \exp[-i\hat{h}(1)\varepsilon/\hbar] | n(X_0) \rangle. \quad (6)$$

In the adiabatic approximation, we pick up the quantum transition only between the states with the same quantum number n ; $\langle n_k | \exp(-i\hat{h}(k)\varepsilon/\hbar) | n_{k-1} \rangle$. Then using the relation $\hat{h}(k)|n_k\rangle = E_n(k)|n_k\rangle$ ($E_n(k)$ is an energy of an adiabatic level n at $X=X_k$), we obtain $T_{nn}(C) = \exp[-i/\hbar \int_0^T E_n(X_t) dt] \langle n(X_0) | n(X_T) \rangle_c$. Here the overlap function $\langle n(X_0) | n(X_T) \rangle_c$ is given as an infinite product:

$$\langle n(X_0) | n(X_T) \rangle_c = \lim_{N \rightarrow \infty} \prod_{k=1}^N \langle n(X_k) | n(X_{k-1}) \rangle, \quad (7)$$

where we adopt a phase convention $|n(X_0)\rangle = |n(X_T)\rangle$. This overlap function naturally involves the history of the excursion in the X -space which is indicated by the suffix C . Each factor $\langle n(X_k) | n(X_{k-1}) \rangle$ in (7) defines a “connection” between two infinitesimally separated points X_{k-1} and X_k , hence Eq. (7) gives a finite connection along circuit C given by a set of division points $\{X_k\}$.⁷⁾ Thus, by using the approximate relation

$$\langle n(X_k) | n(X_{k-1}) \rangle \approx 1 - \langle n | i\partial/\partial X_k | n \rangle \Delta X_k \approx \exp[i\omega],$$

Eq.(7) is written as

$$\langle n(X_0) | n(X_T) \rangle_c = \exp[i\Gamma_n(C)] \quad (8)$$

with

$$\Gamma_n(C) = \oint_C \omega = \oint \langle n | i\partial/\partial X_k | n \rangle dX_k. \quad (9)$$

Equation (9) is essentially the same as the phase obtained by Berry.^{*1} However the present derivation is quite different from Berry's and somewhat similar to Simon's³⁾ which is based upon the holonomy of vector bundle over X -space. Thus we arrive at the effective path integral associated with the adiabatic change of the external dynamical variable X ,

$$K^{\text{eff}}(T) = \sum_n \int \exp \left[\frac{i}{\hbar} (S_n^{ad} + \hbar \Gamma_n(C)) \right] \prod_t d\mu(X_t, P_t), \quad (10)$$

where S_n^{ad} ($\equiv S_0 - \int_0^T E_n(X_t) dt$) is the adiabatic action function. From (10) we get a natural explanation that the phase $\Gamma_n(C)$ appears as a topological action function which is to be added to the usual dynamical action. This is a first consequence of the paper. If we note $\omega = A_i dX_i$ with $A_i = \langle n | i\partial/\partial X_i | n \rangle$, the effective action in (10) can be regarded as the action function for a system in the effective “gauge field” described by the “vector potential” A_i . This result was already obtained by Mead and Truhlar¹⁾ by using the Schrödinger equation which was described by only usual canonical variables (X, P) . However, the present path integral formulation is applicable to more general external systems which are described by non-canonical variables, e.g., spin variables.^{8),9)}

Level-crossing structure revisited: Here we examine a specific model Hamiltonian revealing the topological meaning of the phase Γ ; consider the following internal Hamiltonian:

*1) Here, it is noted that in the present procedure the closed loop C is naturally introduced as a consequence of the trace formula, whereas in Ref. 2) it is presupposed from the outset.

$$\hat{h}(X) = \begin{pmatrix} z & x+iy \\ x-iy & -z \end{pmatrix}, \quad X = (x, y, z). \quad (11)$$

Although this model already has been studied by Berry,²⁾ we treat it in a different way by adopting the "coherent-state" representation,⁹⁾ which may reserve an applicability to more complicated Hamiltonians. Using two-component Pauli spinor we write the eigenvector of (11) as $|\xi\rangle = \cos(\theta/2)|-1/2\rangle + \sin(\theta/2)e^{i\phi}|1/2\rangle = ^t(\sin(\theta/2)e^{i\phi}, \cos(\theta/2)) \equiv ^t(a, b)$, which is given as the $SU(2)$ (spin) coherent-state

$$|\xi\rangle = (1 + |\xi|^2)^{-1/2} \exp[\xi \hat{S}_+] | -1/2 \rangle \quad (12)$$

with $\xi = \tan(\theta/2)e^{i\phi}$. The eigenvalues of (11) are calculated as

$$\lambda_{\pm} = \pm \sqrt{x^2 + y^2 + z^2}, \quad (13)$$

which show a remarkable feature that two levels λ_+ and λ_- cross at $X=0$, namely, the origin $X=0$ becomes a singular-point of cone type. In the following we take the lower level λ_- , for which the corresponding eigenvector is given by

$$a/b = \tan(\theta/2)e^{i\phi} = -(x+iy)/(\sqrt{x^2+y^2+z^2} + z), \quad (14)$$

which yields $\phi = \beta$ and $\theta = -\alpha$, where α and β are the polar angle defined by $x = r \sin \alpha \cdot \cos \beta$, $y = r \sin \alpha \sin \beta$, $z = r \cos \alpha$. $\Gamma(C)$ is thus evaluated as

$$\Gamma(C) = \oint \frac{1}{2i} \frac{(\xi^* \nabla \xi - c.c.)}{1 + |\xi|^2} dX = \frac{1}{2} \oint (1 - \cos \alpha) \nabla \beta dX, \quad (15)$$

which is written as $\oint A dX$, where the vector potential becomes

$$A_x = (\cos \alpha - 1) \frac{\sin \beta}{r \sin \alpha}, \quad A_y = (1 - \cos \alpha) \frac{\cos \beta}{r \sin \alpha} \quad (16)$$

and $A_z = 0$. The striking point is that the negative z -axis (i.e., $\alpha = \pi$) forms a *singular line* on which A diverges.¹⁰⁾ This suggests that the "Dirac pole" is located at the origin as was pointed out by Berry, but the present result gives an explicit form of the specific singular nature. The phase Γ is converted into the surface integral by Stokes' theorem; $\Gamma(C) = \int_S d\omega = \frac{1}{2} \int_S \sin \alpha d\alpha \wedge d\beta$, which is just the solid angle suspended by the closed loop C . Here we note that the present path integral formalism naturally allows the "topological quantization" analogous to the Dirac quantization which is familiar in gauge theory. Consider a sphere S^2 and divide it by C into two hemispheres S and \tilde{S} . We can choose two different gauges such that ω is singular-free on hemispheres S or \tilde{S} , respectively. Then, the topological part of the propagator $K^{\text{eff}}(T)$ can be expressed in two ways according to these two choices, and the consistency condition asserts these two expressions should coincide. Namely, the relation

$$\exp\left[i \oint_C \omega\right] = \exp\left[i \int_S d\omega\right] = \exp\left[-i \int_{\tilde{S}} d\omega\right] \quad (17)$$

should hold, which leads to the quantization condition $\int_S d\omega + \int_{\tilde{S}} d\omega = \int_{S^2} d\omega = 2\pi \times (\text{integer})$.

Finally we give a remark on a generalization to $n \times n$ matrix Hamiltonian; it may be simply achieved by replacing the eigenstate of form (12) by the $SU(n)$ coherent-state⁹⁾

where the parameter space becomes the complex projective space $U(n)/U(n-1) \times U(1)$ the point of which is coordinated by n -dimensional complex vector $\xi = (\xi_1, \dots, \xi_n)$ and the resultant phase yields

$$\Gamma(C) = \frac{i}{2} \oint_{C_{\mu,i}} \sum \left(\frac{\partial \log F}{\partial \xi_\mu} \frac{\partial \xi_\mu}{\partial X_i} - \text{c.c.} \right) dX_i, \quad (18)$$

with $F(\equiv 1 + \xi^\dagger \xi)$ being the overlap function of the coherent state. There may also occur the singularity due to level crossing leading to the general form of the topological quantization.

Semiclassical quantization rule: Now we address a question how one can look at the effect of the topological phase. The most direct way for this is to examine the energy spectra. The energy spectra is rapidly estimated by the semiclassical quantization rule⁽¹⁾ which is derived from the effective propagator (10). Consider the Fourier transform of $K^{\text{eff}}(T)$; $K(E) = i \int_0^\infty K^{\text{eff}}(T) \exp[iET/\hbar] dT$, where we restrict ourselves to a specific adiabatic level n . Firstly the semiclassical limit of $K^{\text{eff}}(T)$ is approximated by the method of stationary phase,

$$K^{\text{sc}}(T) \sim \sum_{P.O.} \exp \left[\frac{i}{\hbar} S^{\text{ad}}(C) + i\Gamma(C) - i\frac{\pi}{2}\alpha(C) \right], \quad (19)$$

where $\alpha(C)$ denotes the so-called Keller-Maslov index and $\sum_{P.O.}$ indicates the sum over periodic orbits. Next, taking the Fourier transform of (19) and evaluating the integral over T by the method of stationary phase, then we get

$$K^{\text{sc}}(E) \sim \sum_{P.O.} \exp \left[\frac{i}{\hbar} W^{\text{ad}}(E) + i\Gamma(C) - i\frac{\pi}{2}\alpha(C) \right], \quad (20)$$

where $W^{\text{ad}}(E) = S^{\text{ad}} + ET$ (action integral) and $T(E)$ is determined by the stationary phase condition $\partial/\partial T(S^{\text{ad}} + ET) = 0$. Here, we restrict ourselves to the case that there appear a finite number of isolated closed orbits for each value of the energy. For this case, a semi-classical quantization condition can be written down explicitly. Namely, taking account of the contribution from the multiple traversals of basic orbits, i.e., putting $W^{\text{ad}} \rightarrow m \cdot W^{\text{ad}}$, $\alpha \rightarrow m \cdot \alpha$ and $\Gamma \rightarrow m \cdot \Gamma$ for m -times traversals and summing over m , $K^{\text{sc}}(E)$ turns out to be

$$K^{\text{sc}}(E) \sim \sum_{P.O.} \exp[i\tilde{W}/\hbar] \cdot \{1 - \exp[i\tilde{W}/\hbar]\}^{-1} \quad (21)$$

with $\tilde{W}(E) = W^{\text{ad}}(E) + \hbar\Gamma(C) - \hbar\alpha/2\pi$. From the pole of (21) we get the formula*

$$W^{\text{ad}}(E) = \oint P dX = \left(n + \frac{\alpha}{4} - \frac{\Gamma}{2\pi} \right) 2\pi\hbar. \quad (22)$$

This gives the energy spectrum for the collective motion including the effect of the topological phase Γ .** Equation (22) is the second main consequence of this paper.

We examine the above formula for a simple case. Consider the internal Hamiltonian

* The more precise form of formula (22) includes stability exponents (see Ref. 11).

** A similar formula has been recently presented by Wilkinson¹²⁾ in the course of investigating the band structure of Bloch electrons in the magnetic field, which is, however, concerned with a rather specific problem. On the other hand, our formula is derived on the general framework of the adiabatic theorem.

of type (11). The level-crossing (singular point) is just located at $x=y=z=0$. Further, we restrict collective periodic orbits to circular motions around the z -axis with a radius ρ , namely collective motions occur in the plane $z=\eta=\text{const}$. Assuming that the collective Hamiltonian has a rotational symmetry, \hat{H}_0 expressed in the polar coordinate is essentially a function of $\hat{p}_\phi = i/\hbar \partial/\partial\phi$ alone. Then, the total Hamiltonian considered is given by

$$\hat{H} = \hat{H}_0(\hat{P}_\phi) + \begin{pmatrix} \eta & \rho e^{i\phi} \\ \rho e^{-i\phi} & -\eta \end{pmatrix}. \quad (23)$$

The B-S phase is $\Gamma_\pm(C) = \mp\pi(1-\eta/e_\pm)$ for the internal upper (lower) state where $e_\pm = \pm\sqrt{\rho^2 + \eta^2}$ denotes an internal energy eigenvalue. From Eq. (22), the semiclassical quantization condition becomes

$$\frac{1}{2\pi} \oint_C P_\phi d\phi = P_\phi(E_\pm) = \left(m - \frac{\Gamma_\pm}{2\pi}\right)\hbar, \quad (24)$$

where $P_\phi(E)$ is defined as an inverse relation of $H_{\text{eff}}(P_\phi)=E$. Semiclassical energy eigenvalues are calculated as

$$\begin{aligned} E_\pm^{(s)} &= H_0\left(P_\phi = m\hbar - \frac{\hbar}{2\pi}\Gamma_\pm\right) + e_\pm \\ &= H_0(P_\phi = m\hbar) + e_\pm \pm \frac{\hbar}{2} \frac{dH_0}{dP_\phi} \Big|_{P_\phi = m\hbar} - \frac{\hbar}{2} \frac{\eta}{e_\pm} \frac{dH_0}{dP_\phi} \Big|_{P_\phi = m\hbar} + O(\hbar^2). \end{aligned} \quad (25)$$

In the present case, \hat{H} can easily be diagonalized and we get exact energy eigenvalues as

$$\begin{aligned} E_\pm &= \frac{1}{2}[H_0(P_\phi = m\hbar) + H_0(P_\phi = (m \pm 1)\hbar) \\ &\quad \pm \sqrt{(H_0(P_\phi = m\hbar) - H_0(P_\phi = (m \pm 1)\hbar))^2 + 4\rho^2}] \\ &= H_0(P_\phi = m\hbar) \pm \sqrt{\rho^2 + \eta^2} \pm \frac{\hbar}{2} \frac{dH_0}{dP_\phi} \Big|_{P_\phi = m\hbar} \mp \frac{\hbar}{2} \frac{\eta}{\sqrt{\rho^2 + \eta^2}} \frac{dH_0}{dP_\phi} \Big|_{P_\phi = m\hbar} + O(\hbar^2). \end{aligned} \quad (26)$$

In Eq. (25), the first and the second term represent the unperturbed eigenvalue of collective Hamiltonian \hat{H}_0 and a conventional adiabatic potential which is constant e_\pm in the present case, respectively. The third and fourth terms come from the B-S phase $\Gamma_\pm(C)$. In comparison with Eq. (26), we can see that this modification of the quantization rule reproduces the exact result up to \hbar -order. Since higher-order terms than \hbar are beyond a semiclassical approximation, eigenvalues (25) can be regarded to be exact within a semiclassical treatment in spite of adiabatic approximation. In addition to the above, Eq. (24) shows that the angular momentum quantization is modified by the existence of the level-crossing. This phenomenon is very analogous to angular momentum quantization in the presence of magnetic flux generated by a solenoid.¹³⁾ Following an analogy with this, Eq. (24) suggests that due to the elimination of internal degrees of freedom the level-crossing produces an “effective spin” for an external dynamical system, which is a sort not to be locally described but only describable by a global aspect of the internal level structure.

Final remarks: The appearance of the topological additive term in the action function

may be regarded as a rather universal phenomenon whenever one deals with interacting dynamical systems in the adiabatic approximation. This viewpoint may shed a light on a wide class of theoretical problems. For example, the "anomaly" in gauge field theory may be naturally understood within the present scheme, which would indeed await a further investigation.

Acknowledgements

The authors would like to thank Dr. T. Hatsuda for a constructive and fruitful discussion. The authors also thank other members of nuclear theory group for their interest. Furthermore, they thank Professor T. Suzuki for his useful comments. One of the authors (S. I.) is indebted to Japan Society for the Promotion of Science for financial support.

References

- 1) C. A. Mead and D. G. Truhlar, *J. Chem. Phys.* **70** (1979), 2284.
- 2) M. V. Berry, *Proc. Roy. Soc. A* **392** (1984), 45.
- 3) B. Simon, *Phys. Rev. Lett.* **51** (1983), 2167.
- 4) D. J. Thouless, M. Kohmoto, M. P. Nightingale and M. den Nijs, *Phys. Rev. Lett.* **49** (1982), 405.
- 5) See, e.g., R. Jackiw, Lecture at Les Houches, July, 1983.
- 6) P. Pechukas, *Phys. Rev.* **181** (1969), 174.
- 7) H. Kuratsuji and T. Hatsuda, in *Proceedings of 13-th International Colloquium in Group Theoretical Method in Physics*, Univ. of Maryland, May 1984, ed. W. W. Zachary (World Scientific), p. 238.
- 8) H. Kuratsuji and T. Suzuki, *J. Math. Phys.* **21** (1980), 472.
- 9) R. Gilmore, *Ann. of Phys.* **74** (1972), 391.
J. R. Klauder, *J. Math. Phys.* **4** (1963), 1055, 1058.
- 10) P. A. M. Dirac, *Proc. Roy. Soc. A* **133** (1931), 60.
- 11) See, e.g., M. C. Gutzwiller, *J. Math. Phys.* **12** (1971), 343.
W. H. Miller, *J. Chem. Phys.* **63** (1975), 996.
- 12) M. Wilkinson, *J. of Phys. A* **17** (1984), 3459.
- 13) F. Wilczek, *Phys. Rev. Lett.* **48** (1982), 1144.

Note added in proof: The non-abelian extension of Berry and Simon's Phase has been recently studied by F. Wilczek and A. Zee (*Phys. Rev. Lett.* **52** (1984), 2111). The authors would like to thank Professor A. Zee for informing of their work.

Adiabatic Effective Lagrangians

John Moody

Department of Computer Science
Yale University
New Haven, CT 06520

Alfred Shapere and Frank Wilczek

School of Natural Sciences
Institute for Advanced Study
Princeton, NJ 08540

ABSTRACT

We discuss the general theory of effective Lagrangians and Hamiltonians in molecular physics. The Born-Oppenheimer effective Lagrangian for the nuclei involves a gauge potential, which may be nonabelian if the electronic energy levels are degenerate. We develop a systematic procedure for finding corrections to the adiabatic approximation, both perturbative and non-perturbative. The former may be incorporated directly into the effective nuclear Lagrangian.

1. The Adiabatic Approximation: General Considerations

Berry's original paper on geometric phases emphasized quantum systems influenced by external parameters [1]. He showed that when these parameters are slowly taken around a closed circuit, the wavefunction of the system may acquire a geometric phase. Although the external parameters were implicitly taken to be classical variables, many interesting applications of the same basic ideas occur in a fully quantum mechanical setting. One can form an effective Born–Oppenheimer Hamiltonian or Lagrangian for the external parameters, that incorporates the effect Berry's phase through a gauge–potential–like term [2] [3] [4]. Upon quantization, the presence of this extra term may lead to significant observable effects, such as shifted quantum numbers and level splittings.

Berry's phase is only the leading correction to the traditional Born–Oppenheimer approximation. Higher-order corrections may also be directly incorporated by adding further terms to the effective Lagrangian [5] [6].

In this paper, we hope to give a relatively unified account of the “modern” Born–Oppenheimer method. We shall discuss both the Hamiltonian and Lagrangian approaches, their relationship, and their apparent inequivalence. Our discussion must necessarily include a description of the procedure for incorporating corrections to the adiabatic approximation, and at the moment, this subject is far from closed. Accordingly, a significant portion of the paper is devoted to a discussion of the adiabatic approximation itself. It covers the Dykhne and Landau–Zener formulas, corrections to them, and geometric phases in the complex plane.

The Born–Oppenheimer approximation first arose in the context of molecular physics [7], but more generally applies whenever a system exhibits two widely separated energy scales. This approximation is often described as a separation of “slow” and “fast” variables; these are just the variables associated with the different energy scales. Quantum mechanically, the separation is made possible by the existence of a large energy gap.

In the original application to molecular physics, the gap involved is the spacing between the electronic energy levels. This gap is typically much larger than the separation between levels associated with vibrations and rotations of the nuclear degrees of freedom that do not involve re-arrangement of electronic orbitals.* Now if we want to describe the spectrum of low-energy excitations of the molecule, *i.e.*, the excitations with energy much less than the electronic energy gap, then we should be able to form a description that involves only the nuclear degrees of freedom. Indeed, at such low energies the electrons have no independent dynamics—they are “enslaved” to

* Often a small finite number of electronic states are actually or approximately degenerate. The formalism appropriate to this case is discussed further below. For now, we assume no degeneracy.

the nuclear degrees of freedom—because only one state is available to them. Therefore, it is possible to describe the low-energy excitations by an effective Lagrangian involving the nuclear degrees of freedom alone, with no explicit reference to the electrons. Of course the value of the numerical parameters appearing in this Lagrangian will depend implicitly upon the electrons.

We find this way of formulating the Born-Oppenheimer idea much more appropriate, and easier to generalize, than the usual formulation in terms of “fast” and “slow” variables. The connection between the two is as follows. Transitions to states separated by a large energy gap require large changes in frequency, and are therefore associated with “fast” variables. Rapid oscillations in time accompany such transitions, and lead to cancellations in processes whose characteristic time scale is much longer—that is, in processes associated with motion of the “slow” variables. Towards the end of this article we shall discuss the relationship between these two approaches more precisely. It is appropriate to mention one conclusion from that discussion now, however: we shall find that quantum variables can only be slow in a very weak sense. For example, in a path integral description the important space-time paths are not differentiable, and the typical velocity is strictly speaking infinite even for so-called “slow” variables. Nevertheless, not being fussy, we shall freely refer to “fast” and “slow” variables throughout this paper.

In quantum field theories containing heavy particles, there is a large gap between the ground state of these heavy particles—*e.g.*, the filled Dirac sea for heavy fermions—and the energy of any excited state. Indeed, to reach an excited state with the same quantum numbers we must in general supply at least the energy to produce a particle-antiparticle pair. Suppose now that the theory contains in addition other fields, describing lighter particles. Then we should be able to describe slow space-time variations of these other fields by an effective Lagrangian that makes no explicit reference to the heavy particles. The usual jargon is to say that we can “integrate out” the heavy particle degrees of freedom. Clearly, the formation of effective Lagrangians in quantum field theory is fully analogous to the corresponding procedure in molecular physics [8]. (But notice that in field theory the heavy particles are the “fast” degrees of freedom!)

In the derivation of effective Lagrangians, we should expect—and will find—that geometric phases occur. This is particularly clear if we think in terms of path integrals. Then along any particular path the slow degrees of freedom can be considered as external parameters governing the state of the fast ones. Therefore, the amplitude for such a path can contain a geometric phase factor of the classic type. Geometric phases of this sort are connected with some of the most subtle and interesting phenomena in quantum field theory, including the occurrence of fractional quantum numbers

and anomalies [9] [10].

2. The Born-Oppenheimer Hamiltonian

We now return to the historical context of the Born-Oppenheimer approximation, to discuss the derivation of effective Hamiltonians and Lagrangians. In molecular physics, it is useful to treat the electronic and nuclear degrees of freedom as fast and slow variables, respectively. This is because the gap between electronic energy levels is typically much larger than the gap between nuclear levels, by a factor of order $(M/m)^{\frac{1}{4}}$ [7]. In the Born-Oppenheimer approximation, one solves for the electronic states in a fixed nuclear background. By the adiabatic theorem, one expects these electronic states to be approximately stationary with respect to the relatively slow motions of the nuclei. We can thus obtain an effective description for the nuclear motion, relative to a fixed electronic orbital, by integrating over electronic coordinates. We shall find that the effective nuclear Lagrangian obtained in this way involves both an ordinary potential term due to the electronic energy levels and a background gauge potential which couples to the nuclear current [2]. This gauge potential takes into account the Berry phase accumulated by the electronic wavefunctions when the nuclear coordinates change adiabatically [3].

The Born-Oppenheimer approximation begins with the full Schrödinger equation

$$(T_{\text{nuc}} + T_{\text{el}} + V)\Psi = E\Psi \quad (2.1)$$

where T_{el} and T_{nuc} are the electronic and nuclear kinetic energy terms, $V(r, R)$ contains the potential and interaction energies of the electrons and nucleons, and r and R are the electronic and nuclear coordinates. The wave function Ψ is separated into nuclear and electronic components Φ_n and ϕ_n as

$$\Psi(r, R) = \sum_n \Phi_n(R) \phi_n(r, R) \quad (2.2)$$

where the subscript n labels the electronic energy eigenstates in a fixed nuclear background. That is, $\phi_n(r, R)$ satisfies the electronic Schrödinger equation at a fixed value of R

$$[T_{\text{el}} + V(r, R)]\phi_n(r, R) = \epsilon_n(R)\phi_n(r, R) \quad (2.3)$$

In terms of the electronic eigenfunctions, the full Schrodinger equation may now be rewritten as

$$\sum_n [T_{\text{nuc}} + \epsilon_n(R)]\Phi_n(R)\phi_n(r, R) = E \sum_n \Phi_n(R)\phi_n(r, R) \quad (2.4)$$

We may now integrate out the electronic degrees of freedom to leave a system of equations for the nuclear wavefunction Ψ alone. Using bracket notation for the normalized electronic eigenstates, we get

$$\sum_n \langle \phi_m | T_{\text{nuc}} \Phi_n | \phi_n \rangle + \epsilon_m(R) \Phi_m = E \Phi_m \quad (2.5)$$

The nuclear kinetic energy operator $T_{\text{nuc}} = -\frac{1}{2M} \nabla^2$ (with $\hbar = 1$) operates on both the nuclear and electronic wavefunctions, $\Phi_n(R)$ and $|\phi_n(r, R)\rangle$. Thus the kinetic energy terms in (2.5) are proportional to

$$\langle \phi_m | \nabla_R^2 \Phi_n | \phi_n \rangle = \sum_k (\delta_{mk} \nabla_R + \langle \phi_m | \nabla_R \phi_k \rangle) (\delta_{kn} \nabla_R + \langle \phi_k | \nabla_R \phi_n \rangle) \Phi_n \quad (2.6)$$

The Born-Oppenheimer approximation applies when the mixing between different electronic levels is small, so that the off-diagonal matrix elements in Eq.(2.6) can be neglected. If, furthermore, the electronic states can be chosen to be real for each R , then $\langle \phi_n | \nabla_R \phi_n \rangle = 0$ and Eq.(2.6) reduces to

$$\left(-\frac{1}{2M} \nabla^2 + \sum_{k \neq n} \frac{1}{2M} \langle \phi_n | \nabla \phi_k \rangle \langle \phi_k | \nabla \phi_n \rangle + \epsilon_n(R) \right) \Phi_n = E \Phi_n \quad (2.7)$$

In this approximation, the nuclei propagate in a background potential

$$\tilde{\epsilon}_n(R) = \epsilon_n(R) + \sum_{k \neq n} \frac{1}{2M} \langle \phi_n | \nabla \phi_k \rangle \langle \phi_k | \nabla \phi_n \rangle$$

The peculiar extra term may be rewritten as follows:

$$\frac{1}{2M} \sum_{k \neq n} \left| \frac{\langle \phi_n | (\nabla H) | \phi_k \rangle}{\epsilon_n - \epsilon_k} \right|^2 \quad (2.8)$$

Hence, when the energy splittings between level n and the other levels are large, this term may be neglected. Berry has pointed out that it is proportional to the trace of the “natural metric” on projective Hilbert space [11].

However, it is not always possible to form a basis of electronic wavefunctions that are everywhere real. Furthermore, corrections to adiabatic evolution will involve mixings of electronic levels. We introduce the “gauge potential” notation

$$A_{mn} \equiv i \langle \phi_m | \nabla_R \phi_n \rangle \quad (2.9)$$

to account for both of these possibilities. Putting together Eqs.(2.4), (2.5), and (2.9), we can write a complete matrix-valued Schrödinger operator for the nuclear wave functions

$$H_{mn}^{\text{eff}} = -\frac{1}{2M} \sum_k (\delta_{mk} \nabla_R - iA_{mk}(R)) \cdot (\delta_{kn} \nabla_R - iA_{kn}(R)) + \delta_{mn} \epsilon_n(R) \quad (2.10)$$

which acts on the vector Φ_n

$$H_{mn}^{\text{eff}} \Phi_n = E \Phi_m \quad (2.11)$$

(The Schrödinger operator is, of course, associated with an effective Hamiltonian after the replacement $-i\nabla_R = p_R$.)

In the Born-Oppenheimer approximation, the effect of the off-diagonal matrix elements A_{mn} which mix different energy levels is ignored. Sections 4 and 5 will be devoted to a justification of this procedure; for now we simply state the result that corrections are indeed suppressed, by a factor depending on the ratio of the typical nuclear and electronic energy splittings. Then for a nondegenerate electronic level, the effective nuclear Schrödinger operator in the Born-Oppenheimer approximation is then simply

$$H_n^{\text{BO}} = -\frac{1}{2M} (\nabla_R - iA_n(R))^2 + \tilde{\epsilon}_n(R) \quad (2.12)$$

where $A_n \equiv A_{nn}$.

Eq. (2.12) looks like the Schrödinger operator of a charged particle in the presence of a background magnetic potential. To further strengthen this analogy, the vector field A_n even transforms like a $U(1)$ gauge potential, as we shall now explain. The phase each of the wavefunctions $|\phi_n(R)\rangle$ is arbitrary, and our description of the dynamics of the nuclei must always respect this arbitrariness. The effect of a redefinition of phases of the electronic wavefunctions $|\phi_n(R)\rangle \rightarrow e^{i\lambda_n(R)}|\phi_n(R)\rangle$, is to rotate the nuclear wavefunctions oppositely

$$\Phi_n(R) \rightarrow e^{-i\lambda_n(R)} \Phi_n(R), \quad (2.13)$$

so that the full wavefunction $\Psi(r, R)$ is preserved. From Eq. (2.9), we see that the gauge potential transforms just as it should:

$$A_n(R) \rightarrow A_n(R) + \nabla_R \lambda_n(R) \quad (2.14)$$

and it is easy to see that the overall effect of the phase redefinition is to leave the nuclear Schrödinger equation invariant (including the term (2.8)).

We conclude that the nuclei behave like charged particles in a magnetic field $B = \nabla \times A_n$. Semiclassically speaking, when the nuclei go around a closed path, the wavefunction will accumulate a geometrical phase proportional to the enclosed magnetic flux. (We will be able to see this more

clearly from the Lagrangian point of view discussed in the following section.) This phase is nothing but Berry's phase in quantum mechanical clothing—the phase that the evolving electron wavefunctions accumulate when their external parameters R are slowly varied has just been passed down to the nuclear wavefunctions.

The degenerate case is slightly more complicated; the evolution will generally involve $U(N)$ rotations among the N degenerate states [12] (provided there are no selection rules forbidding such rotations). In the adiabatic approximation, again restricting attention to a single energy level, we obtain an effective matrix Hamiltonian as in (2.10) with a $U(N)$ gauge potential A_{mn} . The N electronic eigenfunctions may now be regarded as an N -component vector; its “phase” is a $U(N)$ matrix.

For example, the effective nuclear Hamiltonian operator for a molecule with doubly degenerate electronic energy levels (labeled by \uparrow and \downarrow) contains a $U(2)$ gauge potential:

$$H^{\text{BO}} = -\frac{1}{2M} \left\{ \nabla_R - i \begin{pmatrix} A_{\uparrow\uparrow} & A_{\uparrow\downarrow} \\ A_{\downarrow\uparrow} & A_{\downarrow\downarrow} \end{pmatrix} \right\}^2 + \epsilon(R) \quad (2.15)$$

Such a Hamiltonian arises in considering the degenerate Λ -levels of a diatomic molecule [3]. For $\Lambda = \frac{1}{2}$, there is no choice of electronic basis states for which the $U(2)$ gauge potential becomes everywhere diagonal.

To close this section, we remark that there is a much bigger symmetry group that is always present, which mixes states of different energies. The group in question is the unitary symmetry $U(\infty)$ of the electronic Hilbert space \mathcal{H}_R . It is difficult to see where all this symmetry has gone in the nuclear Schrödinger equation Eq. (2.7), because in choosing a decomposition of the total wavefunction in terms of electronic energy eigenfunctions, we have “fixed the gauge” down to a product of $U(N)$ factors (one factor for each N -fold degenerate level). However, there is an alternative formulation in which the full symmetry is manifest, involving a different effective Hamiltonian. We sandwich the time-dependent Schrödinger equation

$$(T_{\text{nuc}} + T_{\text{el}} + V)\Psi = E\Psi \quad (2.16)$$

between a complete set of electronic states (not necessarily energy eigenstates) to obtain a matrix Schrödinger equation analogous to Eq. (2.5)*

$$i(\partial_t - iA_0)\Phi = -\frac{1}{2M}(\nabla_{R_i} - iA_i)^2\Phi + \delta_{mn}\epsilon_n(R) \quad (2.17)$$

where $A_0 \equiv i\langle\phi_m|\dot{\phi}_n\rangle$, with the time derivative referring to the implicit time dependence of $|\phi_n(R(t))\rangle$ (but not on the “dynamical” phase factor

* A similar but not identical nuclear Hamiltonian has been obtained by Zygelman [13].

$\exp i \int \epsilon_n$). Under arbitrary R -dependent (and possibly time-dependent) unitary rotations that preserve (2.2),

$$\begin{aligned}\phi_n(r, R) &\rightarrow \phi_m(r, R) U_{mn}^\dagger(R) \\ \Phi_n(R) &\rightarrow U_{nm}(R) \Phi_m(R)\end{aligned}\quad (2.18)$$

(where $U^\dagger = U^{-1}$), A_0 behaves like the time-component of a gauge field and Eq. (2.17) is fully $U(\infty)$ -covariant. $A_\mu \equiv (A_0, A_i)$ is thus a 4-component $U(\infty)$ gauge potential.

3. Lagrangian Formulation

Often it is more convenient to work with a path-integral description. Phenomenological models are typically easier to formulate in terms of a Lagrangian, where symmetries are manifest. Non-equilibrium and non-perturbative problems, such as calculating tunneling amplitudes, may be easier to solve in the language of path integrals. In addition, as we shall see, it is much easier to incorporate corrections to the adiabatic approximation (which are higher-order in time derivatives) in an effective Lagrangian.

In models of the type we have been considering, effective Lagrangians (typically matrix-valued) for slow degrees of freedom arise naturally when one functional integrates over the fast variables. Generally, functional integration of a matrix-valued integrand requires extra care, to order the operators correctly [14]. But in the adiabatic approximation, if the fast variables are locked into a non-degenerate state, the effective Lagrangian is a scalar, and there is no time ordering to worry about.

As in the previous section, it is convenient to split the Lagrangian into slow and fast parts as follows

$$L = L_{\text{nuc}} + L_{\text{el}} \quad (3.1)$$

with

$$L_{\text{nuc}}(R) = \frac{1}{2} M \dot{R}^2 \quad (3.2)$$

$$L_{\text{el}}(r, R) = \frac{1}{2} m \dot{r}^2 - V(r, R). \quad (3.3)$$

The full time-evolution kernel which connects states at time t_0 to states at time t can be written as a Feynman sum over all paths from configuration (r_0, R_0) to configuration (r, R) :

$$K(r_1, R_1, t_1; r_0, R_0, t_0) = \int_{R_0}^{R_1} \int_{r_0}^{r_1} \mathcal{D}[R] \mathcal{D}[r] \exp i \int_{t_0}^{t_1} dt (L_{\text{nuc}} + L_{\text{el}}(R, r)) \quad (3.4)$$

To form an effective Lagrangian for the nuclei, we now want to integrate (3.4) over the electron coordinates. The first step is to perform a Born-Oppenheimer separation on the molecular wavefunction, as we did in the last section:

$$\Psi(r, R) \equiv \sum_n \Phi_n(R) \phi_n(r, R) \quad (3.5)$$

As before, the electronic eigenstates $\phi_n(r, R)$ are solutions of the electronic Hamiltonian at nuclear configuration R , and the vector-valued nuclear states Φ_n are solutions of the matrix-valued nuclear Schrodinger equation (2.10).

In order to isolate the evolution of the nuclear wavefunctions, we now reorder the general path integral, separating the nuclear and electronic integrations:

$$K = \int_{R_0}^{R_1} \mathcal{D}[R] e^{i \int_{t_0}^{t_1} L_{\text{nuc}}(R) dt} \int_{r_0}^{r_1} \mathcal{D}[r] e^{i \int_{t_0}^{t_1} L_{\text{el}}(r, R) dt} \quad (3.6)$$

With respect to the decomposition of Eq.(3.5), we can express the result of the electronic path integral for a given nuclear path $R(t)$ in terms of an electronic time evolution kernel:

$$\int_{r_0}^{r_1} \mathcal{D}[r] e^{i \int_{t_0}^{t_1} L_{\text{el}}(r, R) dt} \equiv \sum_{mn} \phi_m(r_1, R_1, t_1) K_{mn}^{\text{el}} \phi_n(r_0, R_0, t_0) \quad (3.7)$$

where

$$\begin{aligned} K_{mn}^{\text{el}} &= T \exp \left(-i \int_{t_0}^{t_1} dt \left[\epsilon_m \delta_{mn} + i \langle \phi_m | \dot{\phi}_n \rangle \right] \right) \\ &= T \exp \left(-i \int_{t_0}^{t_1} dt \left[\epsilon_m \delta_{mn} + A_{mn}(R(t)) \cdot \dot{R}(t) \right] \right) \end{aligned} \quad (3.8)$$

is the evolution kernel for the electronic eigenstates. (The notation conveys that we are to take the time ordered exponential of the operator whose mn matrix element is displayed in brackets.) This expression for the kernel comes from integrating the electronic Schrodinger equation for $|\phi_n(R(t))\rangle$

$$\begin{aligned} i \frac{d}{dt} |\phi_n\rangle &= i \frac{\partial}{\partial t} |\phi_n\rangle + i |\dot{\phi}_n\rangle = \epsilon_n |\phi_n\rangle + i \dot{R}(t) \cdot \nabla_R |\phi_n\rangle \\ &= \sum_m \left[\epsilon_n \delta_{mn} + i \langle \phi_m | \nabla_R \phi_n \rangle \cdot \dot{R} \right] |\phi_m\rangle \end{aligned} \quad (3.9)$$

with respect to t . K^{el} just gives the usual dynamical evolution of the electronic energy eigenfunctions with an additional piece coming from the time-dependence of the eigenfunctions through $R(t)$. In the adiabatic limit, the

kernel effectively diagonalizes, and the n th electronic eigenstate obeys the evolution equation

$$|\phi_n(t_1)\rangle = K_{nn}^{\text{el}}(t_1, t_0) |\phi_n(t_0)\rangle = \exp\left(-i \int_{t_0}^{t_1} dt \left[(\epsilon_n + i\langle\phi_n|\dot{\phi}_n\rangle)\right]\right) |\phi_n(t_0)\rangle \quad (3.10)$$

The second term in the exponent is immediately recognized as Berry's phase.

With the electronic motion solved for, the nuclear kernel can now be extracted from the path integral:

$$\Phi_m(R_1) = \left\{ \int_{R_0}^{R_1} \mathcal{D}[R] \text{ T exp } i S[R] \right\}_{mn} \Phi_n(R_0) \quad (3.11)$$

where the exact effective nuclear action is

$$\begin{aligned} S_{mn}^{\text{eff}}[R] &= \int_{t_0}^{t_1} \left\{ \frac{1}{2} M \dot{R}^2 \delta_{mn} - i A_{mn}(R(t)) \cdot \dot{R}(t) - \epsilon_m(R) \delta_{mn} \right\} dt \\ &\equiv \int_{t_0}^{t_1} L_{mn}^{\text{eff}} dt \end{aligned} \quad (3.12)$$

As with the exact effective nuclear Hamiltonian (2.10), the electronic energies $\epsilon_m(R)$ contribute an effective potential for the nuclei, and the velocity-dependent potential term containing A_{mn} modifies the nuclear kinetic energy. In fact, L_{mn}^{eff} can be obtained directly from H_{mn}^{eff} by a Legendre transformation, provided one orders the matrix-valued canonical momenta correctly.

When the electronic levels are nondegenerate, the effective action (3.12) diagonalizes in the Born–Oppenheimer approximation and the time ordering in Eq. (3.11) is unnecessary. For electrons in the n th energy level, the approximate effective action is

$$S_n^{\text{BO}}[R] = \int_{t_0}^{t_1} \left\{ \frac{1}{2} M \dot{R}^2 - i A_n(R(t')) \cdot \dot{R}(t') - \epsilon_n(R) \right\} dt' \quad (3.13)$$

(see also [4]). Curiously, this is *not* the effective action we would have obtained after a Legendre transformation of (2.12), because it does not include the term (2.8). Order by order, the Hamiltonian and Lagrangian formulations of the Born–Oppenheimer approximations are not equivalent.

There are at least two ways one might directly try to incorporate superadiabatic corrections in an effective Lagrangian framework. (The prefix “super‐” indicates that we are looking for corrections to the adiabatic approximation; we are still concerned with the adiabatic regime.) The first is to take, in place of the scalar effective action (3.13), a matrix effective action including a small number of electronic levels, presumably those which are closest in energy to the particular level of interest. The problem with this

scheme is that one still must deal with path-ordered exponentials of matrices, and operator ordering makes the quantization of the nuclear degrees of freedom quite tricky. Another approach is to expand the full time-ordered exponential (3.11) out to some finite degree, and to incorporate this expansion directly into an effective Lagrangian for the n th level, by adding extra terms to (3.13). We shall see how to do this to the lowest super-adiabatic order in the next section.

4. Classical Corrections to Adiabatic Evolution

We now embark on a detailed study of corrections to adiabatic the adiabatic approximation. In this section, our focus will be on corrections to the evolution of electronic wavefunctions with respect to smooth, classical nuclear motions. In Section 5, from the vantage point of the electronic path integral, we briefly discuss why the quantization of the nuclei fundamentally alters the nature and size of these corrections.

The total evolution of a wavefunction which begins in an energy eigenstate is best split into two parts—the amplitude to remain in that eigenstate, and the amplitude to have a transition to another eigenstate. The first part of the problem has been beautifully treated by Berry [15] by means of an iterative procedure. The essential idea is, for a given Hamiltonian $H(t)$, to perform an iterative sequence of time-dependent unitary transformations. At each step, the transformations

$$\begin{aligned}\phi^{(i)}(t) &= U_i(t) \phi^{(i-1)}(t) \\ H_i &= U_i H_{i-1} U_i^\dagger - i U_i \dot{U}_i^\dagger\end{aligned}\tag{4.1}$$

(with $H = H_0$) are supposed to be chosen in such a way that the evolution of the transformed wavefunction with respect to the new Hamiltonian is more adiabatic than the last. The phase evolution of the original energy eigenstate ϕ_n is obtained by evolving

$$\phi^{(i)} = U_i U_{i-1} \cdots U_0 \phi_n\tag{4.2}$$

with respect to H_i in the adiabatic approximation, and then transforming back to the original basis. This scheme is expected to converge rapidly, at least until the i th successive correction becomes comparable in magnitude to the typical amplitude for a transition to another level. At this point, the sequence of iterations begins to diverge—the expansion is asymptotic. (A similar scenario occurs in quantum field theory, where the perturbation expansion is an asymptotic series, which begins to diverge when tunneling processes become important.) Berry's procedure only gives the evolution of the *phase* of ϕ_n ; changes in the magnitude of ϕ_n come from transitions to other levels, which are willfully ignored in this approximation.

In what follows, we shall consider two very different methods for calculating super-adiabatic transition amplitudes. The more straightforward approach is a version of time-dependent perturbation theory (TDPT), that explicitly separates adiabatic from super-adiabatic evolution. The other is essentially non-perturbative and involves analytic continuation into the complex time plane. The perturbative approach will enable us to give a proof of the adiabatic theorem and to find corrections. However, it is not very useful as a calculational tool, and indeed, not always very reliable. On the other hand, the non-perturbative method, embodied in Dykhne's formula and its generalizations, turns out to be quite powerful and accurate. (If some of the following material seems too abstract or technical, the reader may wish to refer to the example beginning with Eq.(4.21) for orientation.)

Adiabatic Perturbation Theory. Our discussion of time-dependent perturbation theory begins with the exact equation for evolution of the electronic wavefunction according to a time-dependent Hamiltonian. From Eqs.(3.8) and (3.9), we have

$$\begin{aligned} |\psi(t)\rangle &= \sum_{mn} |\phi_m(R(t))\rangle U_{mn}(t, -\infty) \langle \phi_n(R(-\infty)) | \psi(R(-\infty))\rangle \\ U(t, -\infty) &\equiv T \exp - i \int_{-\infty}^t dt' [\epsilon_a \delta_{ab} + i \langle \phi_a | \dot{\phi}_b \rangle] \end{aligned} \quad (4.3)$$

where for convenience we have taken $t_1 = -\infty$. In the nondegenerate case, the adiabatic approximation to it is

$$\begin{aligned} |\psi(t)\rangle_{\text{ad}} &= \sum_n |\phi_n(R(t))\rangle \left\{ \exp - i \int_{-\infty}^t dt' [\epsilon_n + i \langle \phi_n | \dot{\phi}_n \rangle] \right\} \\ &\cdot \langle \phi_n(R(-\infty)) | \psi(R(-\infty))\rangle \end{aligned} \quad (4.4)$$

We expect this adiabatic wavefunction to be a better and better approximation to the exact wavefunction, as the time dependence of the parameters $R(t)$ becomes slower and slower. To quantify this, we write

$$R(t, \tau) \equiv R(t/\tau)$$

and let $\psi_\tau(R(t))$ be the solution of the Schrodinger equation when the internal parameters vary as $R(t, \tau)$. Thus, the larger τ is, the more slowly the parameters R are changing. Note that replacing $R(t) \rightarrow R(t/\tau)$ in Eq.(4.3) is equivalent to rescaling time $t \rightarrow \tau t$ everywhere except inside of $R(t)$.

To compare the exact evolution to the adiabatic approximation, it is convenient to use the following general formula for untangling the time-ordered exponential integral of the sum of two operators:

$$\begin{aligned} T \exp \int (M + N) &= T \exp \int M \cdot T \exp \int N' \\ N' &= (T \exp \int M)^{-1} \cdot N \cdot T \exp \int M \end{aligned} \quad (4.5)$$

which is easily proved by differentiating both sides. In our example, we wish to take

$$\begin{aligned} M_{ab} &= -i\epsilon_a \delta_{ab} + \langle \phi_a | \dot{\phi}_a \rangle \delta_{ab} \\ N_{ab} &= \langle \phi_a | \dot{\phi}_b \rangle (1 - \delta_{ab}) \end{aligned} \quad (4.6)$$

so that the true evolution is governed by $M + N$, and the adiabatic evolution by M . By means of these manipulations, the corrections to the adiabatic evolution operator in Eq.(4.3) are isolated into a compact formal expression, that is:

$$\begin{aligned} T \exp \int_{-\infty}^t dt' N'(t') &= T \exp \int_{-\infty}^t dt' \left\{ \langle \phi_p | \dot{\phi}_q \rangle (1 - \delta_{pq}) \right. \\ &\quad \left. \cdot \exp \int_{-\infty}^{t'} dt'' [i(\epsilon_p - \epsilon_q) - \langle \phi_p | \dot{\phi}_p \rangle + \langle \phi_q | \dot{\phi}_q \rangle] \right\} \end{aligned} \quad (4.7)$$

where we have used the fact that M is diagonal to get rid of some of the time ordering.

Our task is now to see how this expression approaches the identity operator as $\tau \rightarrow \infty$. Note first that rescaling the time variable does nothing to the integrals over $\langle \phi_p | \dot{\phi}_p \rangle$ and $\langle \phi_q | \dot{\phi}_q \rangle$, since the scale factors coming from the dt'' in the measure and from the time derivative in the integrand cancel. As we have seen in other contexts, these integrals have a purely geometric character. So, if we rescale the time, the only modification to Eq. (4.7) is to replace $(\epsilon_p - \epsilon_q)$ by $\tau(\epsilon_p - \epsilon_q)$. In other words, the slowness parameter appears only in an oscillatory exponential factor, which we expect will make the total integral very small as $\tau \rightarrow \infty$, à la the Riemann–Lebesgue lemma [16]. To get an idea of just how small, we expand the time-ordered exponential to first order:

$$\begin{aligned} T_{pq} &= T \exp \int_{-\infty}^{\infty} dt' N'(t') \\ &= \delta_{pq} + \int_{-\infty}^{\infty} dt' \left\{ \langle \phi_p | \dot{\phi}_q \rangle (1 - \delta_{pq}) \right. \\ &\quad \left. \cdot \exp \int_{-\infty}^{t'} dt'' [i(\epsilon_p - \epsilon_q) - \langle \phi_p | \dot{\phi}_p \rangle + \langle \phi_q | \dot{\phi}_q \rangle] \right\} + \dots \\ &\equiv \delta_{pq} + \int_{-\infty}^{\infty} dt' \gamma_{pq}(t') \exp i\Delta_{pq}(t') + \dots \end{aligned} \quad (4.8)$$

To simplify this, let us redefine the phases of the wavefunctions $|\phi_n(R(t))\rangle$ for each t so that the wavefunctions are always real along the contour of integration; then $\langle \phi_n | \dot{\phi}_n \rangle = 0$. We now make three crucial assumptions,

that $\epsilon_p \neq \epsilon_q$ for all times, that γ_{pq} and $\epsilon_p - \epsilon_q$ are infinitely differentiable, and that as $t \rightarrow \pm\infty$, all derivatives of γ_{pq} and $\epsilon_p - \epsilon_q$ approach zero with sufficient rapidity. (If $R(t)$ is cyclic, the last assumption is not necessary.) Then integrating by parts, we may ignore all surface terms, and we find for the first-order term in (4.8)

$$T_{pq}^{(1)} = i \int_{-\infty}^{\infty} dt' \frac{\partial}{\partial t} \left(\frac{\gamma_{pq}}{\epsilon_p - \epsilon_q} \right) \exp i \int^{t'} dt'' (\epsilon_p - \epsilon_q) \quad (4.9)$$

from which it follows that

$$|T_{pq}^{(1)}| \leq \left| \int_{-\infty}^{\infty} dt' \frac{\partial}{\partial t} \left(\frac{\gamma_{pq}}{\epsilon_p - \epsilon_q} \right) \right| \quad (4.10)$$

The right-hand side scales like τ^{-1} when we scale $t \rightarrow \tau t$, so $T_{pq}^{(1)}$ goes to zero at least as fast as τ^{-1} . After n repeated integrations by parts, we obtain an expression that vanishes like τ^{-n} . In other words, $T_{pq}^{(1)}$ goes to zero faster than any power of τ . By a similar procedure, one may also show that the same is true for all *off-diagonal* higher-order terms in the expansion (4.8). This completes the proof of the adiabatic theorem.

The diagonal terms in (4.8) need not vanish so fast; indeed, they only vanish like powers of τ . For example, the pp component of the second-order term contains the non-oscillating piece

$$i \sum_{q \neq p} \int_{-\infty}^{\infty} dt' \frac{\gamma_{pq}(t') \gamma_{qp}(t')}{\epsilon_p - \epsilon_q} \quad (4.11)$$

which vanishes only like τ^{-1} . This first-order correction to purely adiabatic evolution may be incorporated directly into the effective Lagrangian (3.13) as a counterterm

$$S_n^{(1)}[R] = \int_{t_0}^{t_1} dt' \left(\sum_{q \neq n} \frac{A_{nq}^i A_{qn}^j}{\epsilon_n - \epsilon_q} \right) \dot{R}_i \dot{R}_j \delta_{mn}. \quad (4.12)$$

It modifies the metric on parameter space—which we have taken to be δ_{ij} —to

$$g_{ij} = \delta_{ij} + \frac{2}{M} \sum_{q \neq n} \frac{A_{nq}^i A_{qn}^j}{\epsilon_n - \epsilon_q}. \quad (4.13)$$

Not surprisingly, Eq.(4.11) is the same expression we would have obtained by expanding the first-order phase approximant from Berry's iteration scheme in powers of τ . When Berry's n th-order phase approximant is rearranged as an expansion in τ , it must agree with the TDPT expansion

through the n th order in τ^{-1} . Both schemes diverge asymptotically. The difference with TDPT is that in principle, we might hope, it may give us information about super-adiabatic transitions.

Unfortunately, the expansion we have presented above is not very useful for actually computing off-diagonal corrections to adiabatic evolution, due to the presence of rapidly oscillating phases and the multiple integrations required to go beyond first order. And even when, say, the first order term in Eq.(4.8) can be evaluated, there may be no guarantee that the second-order term will be smaller—indeed, we shall discuss a specific example below where the higher-order corrections are larger. The moral of the story is that, in the adiabatic regime, transitions between levels are by nature non-perturbative, and attempting to treat them perturbatively is misguided.

Super-adiabatic transitions: Dykhne's formula. Much of what is known about transitions in the adiabatic limit is summarized by an elegant non-perturbative result known as Dykhne's formula [17], relating the amplitude for a transition between two nondegenerate energy levels to the location of their common crossing point in the complex time plane. Suppose that $H(t)$ is a nondegenerate 2×2 Hamiltonian matrix, $E_1(t)$ and $E_2(t)$ are its two instantaneous energy levels, and $E_2 > E_1$ for all real times. If $E_1(t)$ and $E_2(t)$ are extendable into the complex time plane, there will typically be a point t_c where they cross. Dykhne's formula states that the transition probability to go from E_1 to E_2 as t runs from $-\infty$ to $+\infty$ is approximately

$$P_{12} \sim \exp - 2 \operatorname{Im} \int_0^{t_c} (E_2 - E_1) dt \quad (4.14)$$

In general, Dykhne's formula is a good approximation when the crossing points are located far away from the real time axis. This will be true if the energy splittings and/or the typical time scale τ over which $H(t)$ changes are large. The relevant dimensionless expansion parameter is

$$\epsilon \sim \frac{\hbar}{\tau \Delta E} \quad (4.15)$$

and Dykhne's formula states that super-adiabatic transition amplitudes are of order $\mathcal{O}(\exp - \lambda/\epsilon)$ for some positive constant λ . This is the canonical form for non-perturbative corrections to an asymptotic expansion.

It is somewhat ironic that Dykhne's formula for super-adiabatic transitions may be proved by using a version of the adiabatic theorem in the complex time plane [18]. Like the real-time adiabatic theorem, this theorem describes the approximate evolution of the projection of a wavefunction onto an energy eigenstate, along a contour in the complex plane. The idea behind the proof of Dykhne's formula (which we shall only sketch here) is to find an appropriate contour in complex time such that continuation along

the contour connects the two energy levels. For a non-degenerate two-level system, the crossing point will typically be a branch point of square root type for the function $(E_2 - E_1)(t)$. (If $E_2 > E_1$ on the real axis, we will take the branch point in the upper-half plane.) So if a wavefunction is initially in the eigenstate $\phi_1(-\infty)$ with energy $E_1(\infty)$, and if it is evolved along a contour C_{21} which goes over the branch point and across the cut, then one may compute its component in the direction $\phi_2(+\infty)$ (see Fig. 1). This will be related directly to the transition amplitude.

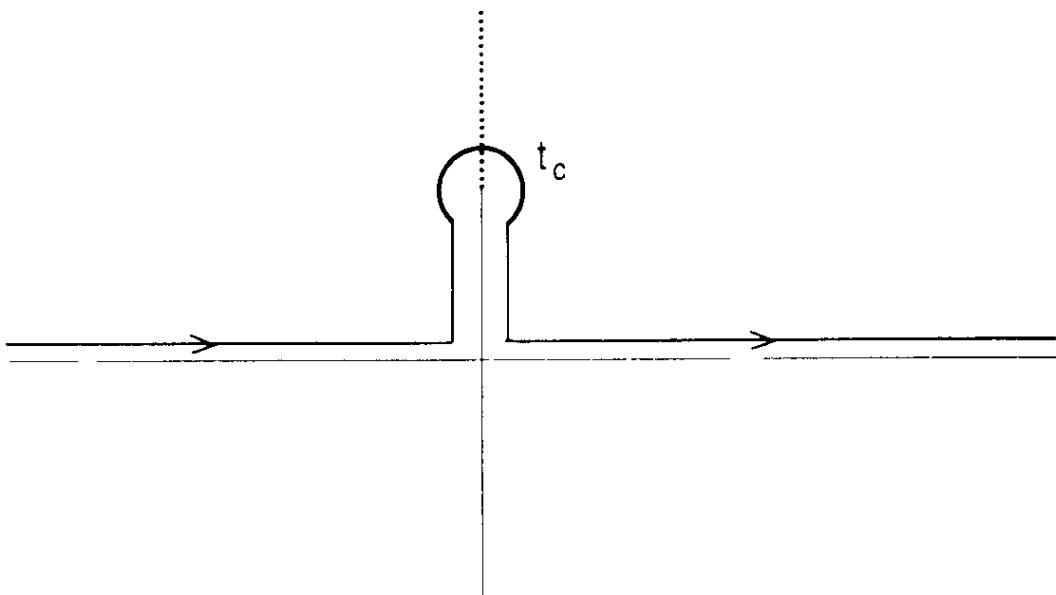


Figure 1. A contour C_{21} that connects two electronic eigenstates with energies E_1 at $t = -\infty$ and E_2 at $t = +\infty$. The function $E_1(t) - E_2(t)$ typically has a square-root branch point in the complex time plane. When the contour passes over the cut, the electron wavefunction crosses over from the E_1 to the E_2 level surface.

To be more explicit, let us choose a smooth basis of eigenstates $\phi_{1,2}(t)$ along the real time axis, and analytically continue our choice into the complex plane. Along the contour C_{21} , we also must have a smooth basis of eigenstates $|\tilde{\phi}_{1,2}\rangle$. Since ϕ_1 and ϕ_2 get interchanged in crossing the cut, it will generally not be possible to take $|\tilde{\phi}_{1,2}\rangle = |\phi_{1,2}(t)\rangle$ everywhere. Instead, our choice along C_{21} will be as follows: to the left of the cut, we take $|\tilde{\phi}_{1,2}(t)\rangle = |\phi_{1,2}(t)\rangle$, but to the right, we take $|\tilde{\phi}_{1,2}(t)\rangle = e^{i\alpha_{2,1}}|\phi_{2,1}(t)\rangle$, where $e^{i\alpha}$ is a phase needed to make the choice of basis continuous across the cut. We shall likewise denote the energy levels along the real time axis by $E_{1,2}(t)$, and along C_{21} by $\tilde{E}_{1,2}(t)$; again because of the cut, $\tilde{E}_{1,2}(+\infty) = E_{2,1}(+\infty)$.

Finally, let $\psi(t)$ be a wavefunction evolving according to $H(t)$, initially in an eigenstate $\psi(-\infty) = \phi_1(-\infty)$. If $H(t)$ is analytic in a strip S , then

$\psi(t)$ will also be analytic and single-valued in S [19]. We may now state the result of the adiabatic theorem for evolution along C_{21} :

$$\langle \tilde{\phi}_1(+\infty) | \psi(+\infty) \rangle \simeq \exp -i \int_{C_{21}} \left[\tilde{E}_1(t) + i \langle \tilde{\phi}_1 | \dot{\tilde{\phi}}_1 \rangle \right] dt \quad (4.16)$$

This is precisely of the same form as Eq.(4.4), sandwiched on the left by $\langle \phi_n |$. Just as Eq.(4.4) said nothing about the “transition” components of ψ , so Eq.(4.16) is silent about the $\tilde{\phi}_2(+\infty) = e^{i\alpha} \phi_1(+\infty)$ component (which in fact is quite large). The big difference here is that the “energies” \tilde{E} need no longer be real. Hence, the norm of (4.16) need not be equal to 1; in fact, in the adiabatic limit, it will be exponentially small.

We now wish to evaluate the transition probability

$$P_{21} \equiv |\langle \phi_2(+\infty) | \psi(+\infty) \rangle|^2 \quad (4.17)$$

The only part of (4.16) contributing to P_{21} comes from the imaginary part of the energy integral. We can put this into a convenient form by deforming the contour so that it follows along the real axis up to $t = 0$, then heads upward to the branch point, then returns to zero and continues along the real axis to infinity. The result is

$$P_{21} = \exp - 2 \operatorname{Im} \int_0^{t_c} (E_2 - E_1) dt \quad (4.18)$$

The transition probability is only part of the story; the phase of the transition amplitude is also of interest. From Eq.(4.16), the phase of the total amplitude is

$$\begin{aligned} & \text{phase}(\langle \phi_2(+\infty) | \psi(+\infty) \rangle) \\ &= e^{i\alpha} \exp - i \int_{-\infty}^0 \left[E_1 + i \langle \phi_1 | \dot{\phi}_1 \rangle \right] dt \\ & \cdot \exp - i \int_0^{t_c} \left[\operatorname{Re}(E_1 - E_2) + i \langle \phi_1 | \dot{\phi}_1 \rangle - i \langle \phi_2 | \dot{\phi}_2 \rangle \right] dt \\ & \cdot \exp - i \int_0^{+\infty} \left[E_2 + i \langle \phi_2 | \dot{\phi}_2 \rangle \right] dt \end{aligned} \quad (4.19)$$

The total phase contains, as usual, both a dynamical and a geometric component. The geometric phase itself splits into two pieces, an adiabatic phase and a *super-adiabatic* phase associated specifically with the tunneling process

$$\exp i\gamma_{1 \rightarrow 2} = \exp \int_0^{t_c} \left(\langle \phi_1 | \dot{\phi}_1 \rangle - \langle \phi_2 | \dot{\phi}_2 \rangle \right) dt \quad (4.20)$$

We shall compute this phase below in the particular case of a Hamiltonian linear in t —the answer will turn out to be independent of any of the parameters appearing in $H(t)$.

The Landau-Zener formula. As an example, we now consider a Hamiltonian $H(t)$ in the vicinity of an avoided crossing, as originally studied by Landau and Zener [20] [21]. We focus upon the two levels whose energies cross, and study the equation governing their mixing in time:

$$i \frac{d\psi}{dt} = H\psi \quad (4.21)$$

$$H(t) = \begin{pmatrix} at & b \\ b & -at \end{pmatrix} = at\sigma_3 + b\sigma_1 \quad (4.22)$$

Here H is the Schrödinger operator in the two-level subspace. We have located the crossing at $t = 0$ and linearized around it, thrown away a possible constant term in the energy, and assumed $b \equiv \langle \psi_p | \psi_q \rangle$ is real; none of these simplifications entails a loss of generality.

The eigenvalues of $H(t)$ are

$$E_{1,2}(t) = \pm \sqrt{b^2 + a^2 t^2} \quad (4.23)$$

and the crossing (a square-root branch point) is located in the upper-half t -plane at $t_c = ib/a$. Hence, according to Eq.(4.18),

$$\begin{aligned} P_{21} &\simeq \exp - 2 \operatorname{Im} \int_0^{t_c} 2\sqrt{b^2 + a^2 t^2} dt \\ &= \exp - \pi \frac{b^2}{a} \end{aligned} \quad (4.24)$$

Dykhne's formula works amazingly well here; in fact, this is the *exact* result obtained by Zener after a much more involved analysis. Incidentally, TDPT is worse than useless here: the first-order term in Eq.(4.8) differs from the correct amplitude by a factor of π .

To find the phase of the transition amplitude requires a little more work. First we need an explicit basis of eigenfunctions: with eigenvalue $-\sqrt{b^2 + a^2 t^2}$ we have

$$\phi_1 = \mathcal{N}_1 \begin{pmatrix} 1 \\ -\frac{at}{b} - \sqrt{1 + \left(\frac{at}{b}\right)^2} \end{pmatrix} \quad (4.25)$$

and with eigenvalue $+\sqrt{b^2 + a^2 t^2}$,

$$\phi_2 = \mathcal{N}_2 \begin{pmatrix} \frac{at}{b} + \sqrt{1 + \left(\frac{at}{b}\right)^2} \\ 1 \end{pmatrix} \quad (4.26)$$

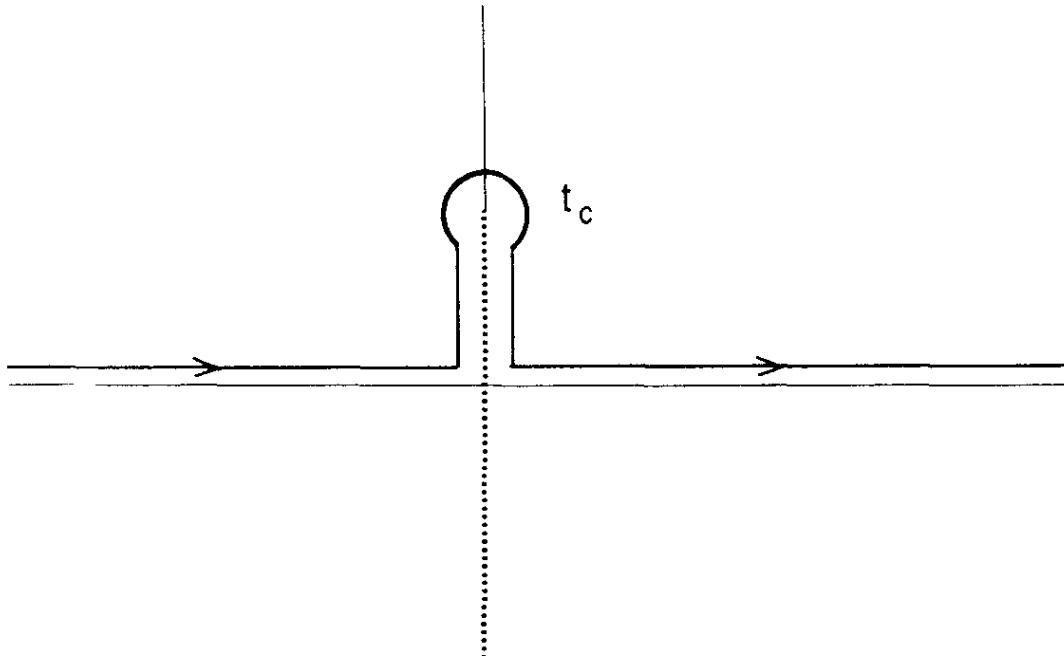


Figure 2. The same contour, from a different vantage point. This figure emphasizes that there is no discontinuity in the evolution of the electron wavefunction along the contour.

For real t , the eigenfunctions may be taken to be real, and the geometric phase receives no contribution from the integration along the real axis. For complex t , we choose the normalization $\mathcal{N}_1 = 1/\sqrt{2}$. We want to evaluate the integral

$$\int_0^{t_c} \left(\langle \phi_1 | \dot{\phi}_1 \rangle - \langle \phi_2 | \dot{\phi}_2 \rangle \right) dt = \int_C \langle \tilde{\phi}_1 | \dot{\tilde{\phi}}_1 \rangle dt \quad (4.27)$$

along the contour shown in Fig. 2. It is important to be sure that the eigenfunctions ϕ_1 and ϕ_2 match up precisely at t_c ; for this to occur, it is necessary to take \mathcal{N}_2 to be $i/\sqrt{2}$. Extending this choice of phase downwards, we find that ϕ_2 is imaginary along the real t -axis. The actual computation of the geometric phase (4.20) is straightforward; the result is

$$\gamma_{1 \rightarrow 2} = -i \frac{\pi}{2} \quad (4.28)$$

for a total phase of $-i$. This phase is precisely what is needed to cancel the i picked up in matching the wavefunctions at the branch point; it makes the non-dynamical part of the wavefunction real for all real times.

It is curious that the geometric phase we just computed is completely independent of a and b . In fact, we can argue that this sort of phase will arise quite generally, for Hamiltonians that are real on the real time axis. Whenever the energies cross at t_c , they will also cross at the conjugate point t_c^* . Let us join these two branch points by a cut drawn contours above and below the cut, in such away that the images of the two contours are complex

conjugate (see Fig. 3). It is easy to see that the total phases obtained by integrating along either contour must also be conjugate, and since the wavefunction must be single-valued, the two phases must be equal to ± 1 . Furthermore, if the contour is chosen so that the total dynamical phase vanishes, then the geometric part of the wavefunction must be real, as we found in our example above. (This argument of course does not apply to complex Hamiltonians, and in general we can obtain complex geometric phases for such processes.)

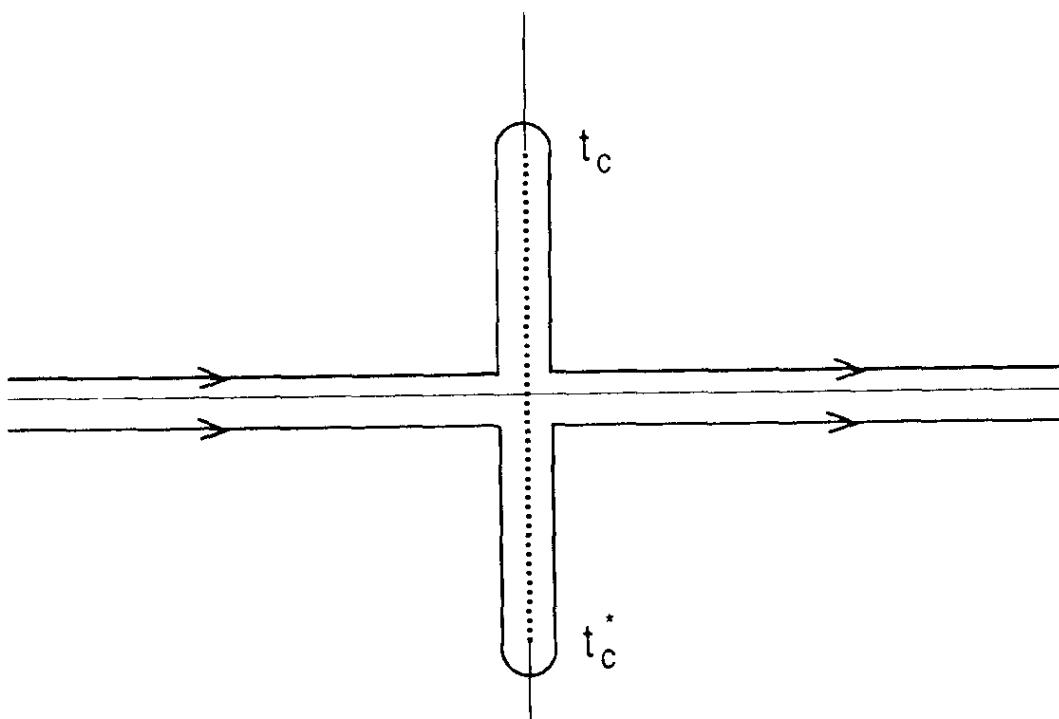


Figure 3. Real Hamiltonians will have crossings at conjugate points. The results of evolution along the lower and upper contours are conjugate, and equal.

Corrections to tunneling. We would briefly like to discuss corrections to Eqs.(4.18) and (4.19). The simplest types of corrections will come from crossing points farther from the real time axis. These will be of the same form as before, with a sum over all of the other crossing points.

Less straightforward are corrections to adiabatic evolution along a contour, modifying a super-adiabatic tunneling amplitude. To handle these, one may follow a similar prescription to the above perturbative expansion for diagonal corrections to adiabatic evolution. As in Eq.(4.7), one may separate off the adiabatic part and expand the residual path-ordered exponential; the result will be a series expansion in powers of τ . The zeroth-order term, as we have seen, is a pure geometric phase, but higher-order terms will correct both the phase and norm of the transition amplitude.

5. Quantum Corrections

In the preceding section, we discussed perturbative and non-perturbative corrections to classical adiabatic evolution. We found that corrections to the phase evolution of an energy eigenstate could be put into the form of an asymptotic expansion in powers of τ^{-1} , and that transitions could be calculated by analytically continuing a smooth nuclear path into the complex time plane.

Now we would like to understand how these conclusions are affected when the nuclear degrees of freedom are quantized. In the previous section, we assumed that the nuclear motions were infinitely differentiable, in order to derive our expansion (which involved performing successive integrations by parts) and to prove that tunneling corrections to adiabatic evolution vanish faster than any power of τ^{-1} . The delicate argument breaks down when we try to integrate over nuclear paths in quantum mechanics, because a typical path in the measure is generically *nowhere* differentiable [22]. Needless to say, for such a path, our analytic continuation method for calculating tunneling amplitudes does not apply. Furthermore, in a quantum mechanical context, there is no reason to expect the individual terms in our perturbation series to converge. Indeed, each successive term added to the effective Lagrangian, being higher-order in time derivatives than any of the preceding terms, represents a singular perturbation and diverges for a typical path. How is all this consistent with the successful use of the Born–Oppenheimer approximation in quantum mechanics? Two questions need to be asked: Do we still have a useful perturbative expansion in powers of \dot{R} (and higher time derivatives) for the evolution of an eigenstate? Are tunneling processes still exponentially suppressed?

In general, this will only be true if we take matrix elements between particularly nice states. The point is the following. In passing from our effective Lagrangian to a Hamiltonian, the powers of \dot{R} are converted into powers of the momentum p . (The procedure here is to treat the higher-derivative terms as perturbations, ignoring their effect on the canonical momenta, and to re-express them in terms of $p = M\dot{R}$.) Now p is an *unbounded* operator, and so our expansion “typically” diverges. However, for suitable initial and final states, p may have small matrix elements, and in that case our adiabatic expansion is useful. Similar remarks apply to semiclassical expansions around *smooth* tunneling paths.

To conclude, the validity of the adiabatic approximation in situations where the external parameters are themselves quantized is far from obvious, and should be studied on a case-by-case basis. Nevertheless, in many useful cases, the corrections are expected to be small.

We would like to thank William Bialek, Joanne Cohn, David Eliezer, Paweł Mazur, Hirosi Ooguri, and Bernard Zygelmans for useful discussions.

A significant part of this research was done at the Institute for Theoretical Physics of the University of California, Santa Barbara. This research was supported in part by the National Science Foundation under Grant No. PHY82-17853, supplemented by funds from the National Aeronautics and Space Administration. A.S. was also supported in part by the Department of Energy under Grant. No. DE-AC02-76ERO-2220

References

- [1] M.V. Berry, "Quantal phase factors accompanying adiabatic changes," *Proc. R. Soc. Lond. A.* **392** (1984) 45-57.
- [2] C.A. Mead and D.G. Truhlar, "On the determination of Born-Oppenheimer nuclear motion wave functions including complications due to conical intersections and identical nuclei," *J. Chem. Phys.* **70**, (1979) 2284-96.
- [3] J. Moody, A. Shapere, and F. Wilczek, "Realizations of magnetic-monopole gauge fields: Diatoms and spin precession," *Phys. Rev. Lett.* **56**, (1986) 893.
- [4] H. Kuratsuji and S. Iida, "Effective action for adiabatic processes," *Prog. Th. Phys.* **74**(1985) 439-445
A. Bulgac, "Effective action for nonadiabatic processes," *Phys. Rev. A* **37**, (1988) 4084.
- [5] R. Jackiw, "Three Elaborations on Berry's Connection, Curvature and Phase." *Int. J. Mod. Phys. A***3** (1988) 285-297.
- [6] Ph. de Sousa Gerbert, "A systematic expansion of the adiabatic phase," MIT preprint no. CTP-1537, submitted to Nuclear Physics B.
- [7] A. Messiah, *Quantum Mechanics*, vol.2 (Amsterdam: North-Holland).
- [8] F. Wilczek, lectures delivered at the Theoretical Advanced Study Institute in High-Energy Physics, Univ. of Michigan, 1984.
- [9] Chapters 5 and 7, this volume.
- [10] M. Stone, "Born-Oppenheimer approximation and the origin of Wess-Zumino terms: Some quantum mechanical examples," *Phys. Rev. D* **33** (1986) 1191;
P. Nelson and L. Alvarez-Gaume, "Hamiltonian interpretation of anomalies," *Commun. Math. Phys.* **99**, (1985) 103-114.
- [11] M.V. Berry, "The quantum phase, five years later," Chapter 1, this volume.
- [12] F. Wilczek and A. Zee, "Appearance of gauge structure in simple dynamical systems," *Phys. Rev. Lett.* **52** (1984) 2111.
- [13] B. Zygelman, "Appearance of gauge potentials in atomic collision physics," *Phys. Lett.* **125A**, (1987) 476.
- [14] L. Schulman, *Techniques and Applications of Path Integration* (New York: Wiley, 1981).

-
- [15] M.V. Berry, "Quantum phase corrections from adiabatic iteration," *Proc. R. Soc. Lond. A* **414** (1987) 31–46.
 - [16] M. Reed and B. Simon, *Functional Analysis* (New York: Academic Press, 1972).
 - [17] A.M. Dykhne, "Adiabatic perturbation of discrete spectrum states," *JETP* **14** (1962) 941.
 - [18] J-T. Hwang and P. Pechukas, "The adiabatic theorem in the complex plane and the semiclassical calculation af nonadiabatic transition amplitudes," *J.Chem.Phys.* **67** (1977) 4640-4653.
 - [19] E.A. Coddington and N. Levinson, *Theory of Ordinary Differential Equations* (New York: McGraw-Hill, 1955).
 - [20] A. Zener, "Non-adiabatic crossing of energy levels," *Proc.Roy.Soc.A* **137**, (1932) 696.
 - [21] L.D. Landau and E.M. Lifshitz, *Quantum Mechanics* (Oxford: Pergamon, 1977) section 52.
 - [22] L.F. Abbott and M.B. Wise, "Dimension of a quantum-mechanical path," *Am.J.Phys.* **49(1)** (1981) 37–39.

5

Fractional Statistics

Consider a particle of mass m and charge q moving on a circular ring of radius R . Suppose that the ring is threaded by a thin solenoid carrying flux ϕ . Then the Lagrangian for the particle is¹

$$L = \frac{1}{2}mR^2\dot{\theta}^2 + \frac{q\Phi}{2\pi}\dot{\theta} \quad (5.1)$$

From this we may read off the canonical momentum

$$p_\theta = mR^2\dot{\theta} + \frac{q\Phi}{2\pi} \quad (5.2)$$

and the Hamiltonian

$$H = \frac{1}{2mR^2} \left(p_\theta - \frac{q\Phi}{2\pi} \right)^2 \quad (5.3)$$

The eigenfunctions are

$$\psi_n = e^{in\theta} \quad (5.4)$$

with energies

$$E_n = \frac{1}{2mR^2} \left(n - \frac{q\Phi}{2\pi} \right)^2 \quad (5.5)$$

The energy is easily interpreted as half the square of the kinetic angular momentum, divided by the moment of inertia. Notice that the allowed kinetic angular momenta are spaced by integers—that is, by whole multiples of Planck's constant—but are uniformly displaced from integers by $q\Phi/2\pi$, which might be any real number.

Now all the Lagrangians (5.1), for different values of Φ , yield the same classical equations of motion. Indeed, the term by which they differ is a total derivative, and cannot affect the classical variational principle. As we have seen, however, they lead to different quantum theories. This is not difficult to understand, from the point of view of canonical quantization. To canonically quantize the theory, we need not only to derive the equations of motion but also to impose commutation relations. The commutation relations are

imposed between the coordinates and their conjugate momenta, and their content is changed if the definition of these momenta is changed, whether or not the equations of motion are affected. In the simple quantization problem above, we have seen how just such a modification can affect the physical consequences of a given classical Lagrangian.

The situation is slightly puzzling, however, if we consider it from the point of view of path integral, as opposed to canonical, quantization. In path integral quantization, transition amplitudes are calculated by adding the contributions of all possible paths, each weighted by the exponential of the classical action along the path. Thus it would seem that the transition amplitudes, and thereby the entire physical content of the theory, should be determined by the classical action. On the other hand, we have just seen that different Lagrangians, equivalent at the classical level, can lead to different quantum theories. How did this occur? And could we have foreseen this possibility in advance, knowing only the degrees of freedom and the classical limit? The latter question, as we shall see, is of some practical importance, since in proposing effective Lagrangians for complicated systems we often have a clear idea what the important degrees of freedom are, and expectations for their behavior in the classical limit, but a much more shadowy idea of anything more subtle. In such situations, we would like to know whether quantization is ambiguous, and what further considerations are necessary to resolve the ambiguity.

Any reader who has followed us to this point, or kept the title of this book in mind, will naturally suspect that the geometric or Berry phase plays a key role in this sort of problem. Indeed, its quasi-kinematic character, and the fact that it forms the first (order \hbar^0) correction to purely classical (order \hbar^{-1}) behavior, strongly hint at this role.

Let us return from these generalities to the case at hand. Since we have the explicit Lagrangian (5.1) in front of us, we can easily get to the root of the problem. It is that the classically ignorable $\dot{\theta}$ term is not ignorable in the quantum theory. What does this term do, path by path? It suffices to compare paths that begin at a common position θ_1 at time t_1 and end at a definite position θ_2 at time t_2 , since only such paths can interfere. It is easy to see that the effect of the $\dot{\theta}$ term for such paths is to weight their relative contribution by

$$\exp \left[i \frac{q\Phi}{2\pi} \left(\int_{\text{path 1}} \dot{\theta} dt - \int_{\text{path 2}} \dot{\theta}_2 dt \right) \right] \equiv \exp i \frac{q\Phi}{2\pi} \delta\theta \quad (5.6)$$

Now $\delta\theta$ is a multiple of 2π for paths with common endpoints (and only such paths can interfere). Loosely speaking, $\delta\theta/2\pi$ measures the difference between the number of times the first and second paths winds around the ring, respectively. More precisely, it is the number of times the composite path we get by following the first path from t_1 to t_2 and the inverse of the second path back from t_2 to t_1 winds around the ring.

Working backwards, we can now see very clearly why there was an ambiguity in formulating the path integral. By the argument given immediately above, we can concentrate on closed paths—that is, paths that begin and end at the same point. Such paths fall into distinct, disconnected classes, labeled by the winding number. Because the winding number is an integer for any path, continuous changes in the path cannot alter its value at all. Now, the crucial observation is that the classical Lagrangian cannot guide us in choosing how to weight the relative contributions of disconnected classes of paths. For the classical equations of motion follow from a variational principle that involves only comparisons among infinitesimally nearby paths, and cannot give guidance in comparing disconnected classes of paths. And so, if the closed paths in configuration space fall into disconnected classes, then there is an ambiguity in quantization. In standard mathematical language, we would say that there is a possibility for a single classical Lagrangian to lead to various quantum theories, when the first homotopy group of the configuration space is non-trivial, or in other words when the configuration space is not simply connected.

$$e^{iS_{1+2}} = e^{iS_1} \cdot e^{iS_2}$$

Figure 5.1 The amplitude for the composition of two paths is the product of the amplitudes for each path separately.

There is an important principle constraining the choice of relative weightings between different classes of closed paths. It has to do with the “group” aspect of the homotopy group.² According to basic principles of quantum mechanics, the amplitude for the composition of two paths must be the product of the amplitude for each path separately (see Figure 5.1). This rule is manifestly obeyed by the usual assignment, that weights with the exponential of the integral of the classical action—since the integrals add, the exponentials multiply. We must make sure that any additional weighting factors obey the rule separately. Now composing paths is precisely the operation that defines the homotopy group. The elements of this group are

the disconnected classes of paths, and the product of two such classes is obtained by composing representative paths from each; it is the class of the composed path. Thus if we are assigning extra numerical factors α_π to the paths, we must demand that they obey the rule

$$\alpha_{\pi_1 \circ \pi_2} = \alpha_{\pi_1} \cdot \alpha_{\pi_2} \quad (5.7)$$

or in other words that the factors form a (one-dimensional) representation of the homotopy group. Higher dimensional representations may also have a sensible interpretation; the states have then a non-classical internal degree of freedom, that locates them in the representation space. Some issues that arise in the quantization of theories on non-simply connected spaces are addressed in recent work of Balachandran.³

The phase factors we have been discussing certainly deserve to be called geometric—in fact, being topological, they go one step further towards pure kinematics. Are they related to the geometric phases discussed in previous chapters? In fact, the relationship is quite close. Let us illustrate this on our example. The geometric phase between nearby position eigenstates $|\theta\rangle$ and $|\theta + \Delta\theta\rangle$ is, formally

$$\exp i \frac{q\Phi}{2\pi} \Delta\theta \quad (5.8)$$

This phase is, of course, locally integrable, but globally it is not—when we come around adiabatically by 2π , we have accumulated a phase of $e^{iq\Phi}$. Thus the geometric phase around a homotopically non-trivial path in configuration space parametrizes the ambiguity in quantization. The relationship demonstrated here in a simple example is much more general, as we shall see when we come to discuss anomalies and Wess-Zumino terms in Chapter 7.

An interesting implication of the above is that the wave function in general cannot be defined on ordinary configuration space. For as we have seen, the position eigenstates for θ and $\theta + 2\pi$ are related by a phase factor. Properly then, we should define the wave function with θ running from $-\infty$ to $+\infty$, with the boundary condition

$$\psi(\theta + 2\pi) = e^{iq\Phi} \psi(\theta) \quad (5.9)$$

In general, the wave function will live on the universal covering space of configuration space (or, if there are internal degrees of freedom, a vector bundle over the covering space) and will obey boundary conditions relating points that project to the same point in configuration space. Actually, we have seen examples of this arise before, abstractly in the case of the monopole bundle (introduction to Chapter 3) and concretely in the context of molecular physics (Chapter 4), where wave functions for half-integral orbital angular momentum appeared. Such angular wave functions cannot be realized as a vector bundle of the familiar, “trivial” kind (concretely, as an array of $2j+1$

functions on the sphere, the vector spherical harmonics) but live comfortably on a twisted bundle.

A notable application of these ideas arises in the quantum mechanics of identical particles. The configuration space for N identical particles in d -dimensional space is

$$\overbrace{\mathbf{R}^d \times \cdots \times \mathbf{R}^d}^{N \text{ copies}} / S_N \quad (5.10)$$

The notation indicates that points differing by a permutation of positions are to be identified. Let us assume also that the amplitude for particles to collide vanishes. (This is found to be true *a posteriori* for anything but bosons, because there is an effective centrifugal barrier. For bosons, the topology of configuration space collapses, but that doesn't matter since the paths are equally weighted anyway.) Then we arrive at the configuration space

$$(\mathbf{R}^d \times \cdots \times \mathbf{R}^d - \Delta) / S_N \quad (5.11)$$

where Δ is the space of configurations where two or more particles occupy the same point in \mathbf{R}^d .

To illustrate, let us consider the case $N = 2$ in more detail. In general, for $d > 2$ the homotopy group of the configuration space for N identical particles is simply the group of permutations. For $d = 2$, however, it is a much more interesting group, the so-called braid group. As we have seen, new possibilities for quantization are associated with representations of the homotopy groups. The one-dimensional representations of the permutation group are of course well known; they are the trivial representation and the "sign" representation, that assigns a factor $+1$ or -1 to even and odd permutations respectively. These two possibilities correspond to quantizing the identical particles as bosons or fermions, respectively. The one-dimensional representations of the braid group are analyzed in the enclosed paper by Wu [5.3]. It turns out that they are parametrized by a single parameter, a complex number of unit modulus, commonly written $e^{i\theta}$. When this parameter is $+1$ we have bosons, when it is -1 we have fermions; the peculiar feature of two spatial dimensions is the possibility of continuous interpolation between these familiar cases.

Higher dimensional unitary representations of the homotopy groups also exist. For $d > 2$, particles quantized in this way are said to obey parastatistics. The higher-dimensional unitary representations of the braid group are, it seems, not completely classified. They arise in conformal field theories, and are a very active topic of recent research.⁴

There is a very simple and attractive way of realizing fractional statistics, closely related to the example discussed above. Indeed the characteristic feature of the problem we encountered in quantizing the charged particle on a flux-enclosing ring was the possibility of interactions which generated

phase factors for non-contractible paths. To implement fractional statistics, such phase factors are precisely what we need. All this suggests that fractional statistics—or, in the jargon, non-integrable holonomy factors on configuration space—can be realized dynamically, by a suitable arrangement of fictitious charge and flux. Such a construction has been performed in great generality, both for quantum mechanics and quantum field theory, in the enclosed paper of Arovas, Schrieffer, Wilczek and Zee [5.4], and is discussed further in the article by Arovas in the following chapter [6.2].

Finally, let us mention several occurrences of fractional statistics in the natural world.

Effectively two-dimensional tubes of magnetic flux play an important role in the theory of type II superconductors. Indeed, magnetic fields penetrate type II superconductors only in such tubes. The amount of flux carried by a tube is quantized in units of $\Phi_0 = 2\pi/2e$. This is easily understood in terms of our particle on a ring example: the lowest value of the energy can only be achieved for integral values of $q\Phi/2\pi$ in equation (5.5). For a superconductor we are concerned with the wave function for Cooper pairs, so that $2e$ is the relevant charge. Also, the quantization is sharp (for a macroscopic sample) because a very large number of particles are involved in the condensate. Thus the only acceptable way of accomodating magnetic flux, that does not involve energy proportional to the volume of the sample, is to confine it in small flux tubes carrying integral multiples of Φ_0 . Now if we consider a single electron—not screened by the condensate—orbiting around a flux tube carrying flux ϕ_0 , we see that this composite provides one realization of our “statistical” interaction. Both the flux tube, and the flux tube plus an electron, are bosons. Closely related (but more elaborate and speculative) phenomena occur for particles orbiting string solutions in unified models of particle physics.⁵

The non-linear sigma model, which assigns to each point in space a direction, is used in the description of the low-energy excitations of magnetically ordered systems. In two spatial dimensions, its configuration space is not simply connected, a fact closely connected with the existence of the topological construct known as the Hopf invariant. As a result, the quantization of this classical model of magnetism is ambiguous. This situation is analyzed in the enclosed paper of Wilczek and Zee [5.2]. The correct quantization must be determined by appealing to the microscopic model underlying the effective sigma-model. This model has been a subject of recent research, given some urgency by the importance of two-dimensional magnetism in materials supporting high temperature superconductivity.

Finally, perhaps the most intriguing application of fractional statistics to date concerns the quantized Hall effect. This application is so interesting in itself, and such a nice example of the use of the geometrical phase, that we have devoted the following Chapter exclusively to it.

-
- [1] F. Wilczek, *Phys. Rev. Lett.* **48** (1982) 1144.
 - [2] N.D. Mermin, "The topological theory of defects in ordered media," *Rev. Mod. Phys.* **51** (1979) 591–648.
Charles Nash and Siddhartha Sen, *Topology and Geometry for Physicists*, London: Academic Press, 1983.
 - [3] A.P. Balachandran, "Topological Aspects of Quantum Gravity," Syracuse Univ. preprint no. SU-4228-373.
 - [4] Cumrun Vafa, "Toward Classification of Conformal Theories," *Phys. Lett.* **206B** (1988) 421.
Edward Witten, "Quantum Field Theory and the Jones Polynomial," Inst. for Advanced Study preprint no. IASSNS-HEP-88/33.
Gregory Moore and Nathan Seiberg, "Classical and Quantum Conformal Field Theory," IAS preprint no. IASSNS-HEP-88/35.
 - [5] Mark Alford and Frank Wilczek, "Aharonov–Bohm Interaction of Cosmic Strings with Matter," Harvard Univ. preprint no. HUTP-88/A047, to appear in *Physical Review Letters*.

Linking Numbers, Spin, and Statistics of Solitons

Frank Wilczek

Institute for Theoretical Physics, University of California, Santa Barbara, California 93106

and

A. Zee

*Institute for Theoretical Physics, University of California, Santa Barbara, California 93106, and
Department of Physics, (a) University of Washington, Seattle, Washington 98195*

(Received 24 October 1983)

The spin and statistics of solitons in the $(2+1)$ - and $(3+1)$ -dimensional nonlinear σ models is considered. For the $(2+1)$ -dimensional case, there is the possibility of fractional spin and exotic statistics; for $3+1$ dimensions the usual spin-statistics relation is demonstrated. The linking-number interpretation of the Hopf invariant and the use of suspension considerably simplify the analysis.

PACS numbers: 02.40.+m, 11.10.Lm, 11.30.-j, 75.70.-i

The existence of solitons in the $(2+1)$ -dimensional $O(3)$ nonlinear σ model, as first discussed by Belavin and Polyakov,¹ is implied by the homotopy $\pi_1(s_2) = \mathbb{Z}$. The model is described by the functional

$$E = \int d^2x (1/2f)(\partial_i n^a)^2, \quad i=1, 2; \quad a=1, 2, 3, \quad (1)$$

giving the energy of a static configuration specified by $n^a(\vec{x})$. The "order parameter" n^a is a three-dimensional unit vector: $n^a n^a = 1$. If we describe the ground state by $\hat{n}(\vec{x}) = (0, 0, 1)$, then the basic soliton is described by

$$\hat{n}(\vec{x}) = (\hat{x} \sin f, \cos f). \quad (2)$$

Here $\hat{x} = \vec{x}/|\vec{x}|$ denotes the two-dimensional unit radial vector and $f(r = |\vec{x}|)$ is a function varying smoothly and monotonically from $f(0) = \pi$ to $f(\infty) = 0$ as r increases. We refer to such a topological configuration as a skyrmion.² The topological current in this model is

$$J^\mu = (1/8\pi) \epsilon^{\mu\nu\lambda} \epsilon^{abc} n^a \partial_\nu n^b \partial_\lambda n^c. \quad (3)$$

The space-time indices μ, ν, \dots run over 0, 1, 2. One easily verifies the conservation of this. The topological charge of this current,

$$Q = \int d^2x J^0 \\ = (1/8\pi) \int d^2x \epsilon^{ij} \epsilon^{abc} n^a \partial_i n^b \partial_j n^c, \quad (4)$$

clearly describes the homotopy of the mapping $s_2 \rightarrow s_2$ for \hat{n} satisfying the boundary condition $\hat{n}(\vec{x} = \infty) = \text{const}$. The skyrmion has $Q = 1$. By using Bogomolny's inequality, one can solve exactly the problem of minimizing the energy functional for a given Q . Finally, we mention that this model provides a phenomenological description of Heisenberg ferromagnets in a two-dimensional system and thus the phenomena exhibited in this mod-

el may conceivably be accessible experimentally.

In this paper we point out that the skyrmion may possess fractional angular momentum and obey peculiar quantum statistics. One of us had previously proposed⁴⁻⁶ the possibility of fractional angular momentum and of statistics which are neither Bose-Einstein nor Fermi-Dirac. As we will see, the $(2+1)$ -dimensional $O(3)$ nonlinear σ model provides an amusing and explicit field-theoretic realization of these ideas. Our discussion is also related to several other field-theoretic phenomena discovered in recent years.

The relevant mathematics which allows skyrmions to have these peculiar properties is the homotopy $\pi_3(s_2) = \mathbb{Z}$ [which is perhaps somewhat less obvious than the homotopy $\pi_2(s_2) = \mathbb{Z}$ responsible for the skyrmion's existence]. It is easy to exhibit the basic Hopf map of $s_3 \rightarrow s_2$. In fact, physicists should be familiar with this fact from elementary discussions of the Pauli matrices σ^a . Define $n^a = z^\dagger \sigma^a z$ where

$$z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

is a complex two-component spinor with the constraint $|z_1|^2 + |z_2|^2 = 1$. Notice that the $U(1)$ transformation $z \rightarrow e^{i\theta} z$ leaves n invariant and so the inverse image of any point on s_2 is a circle on s_3 . What is not so obvious is the construction of a Hopf invariant to describe $\pi_3(s_2)$ [just as Q describes $\pi_2(s_2)$]. This will be explained below.

Let us first address the physical question of the spin of the skyrmion. To determine the spin, we rotate the skyrmion adiabatically through 2π over a long time period T . According to Feynman, at the end of this rotation the wave function acquires a phase factor e^{iS} where S is the action corre-

sponding to the adiabatic rotation. The angular momentum J of the skyrmion is given by $e^{is} = e^{i2\pi J}$.

Now, if S has simply the standard form [cf. Eq. (1)]

$$S_0 = \int d^3x (1/2f) (\partial_\mu n^\alpha)^2, \quad (5)$$

then it is easy to see that S_0 is of order $1/T - 0$ as $T \rightarrow \infty$. The skyrmion has $J = 0$. However, we have not taken into account the possibility of including in S a topological term. This possibility is by now familiar from the discussion of the θ vacua⁷ in quantum chromodynamics, from studies of three-dimensional Yang-Mills theory and gravity,⁸ and from recent work of Witten⁹ on strongly interacting skyrmions¹⁰ (based on earlier work of Wess and Zumino¹¹). In general, we can have $S = S_0 + \theta H$ where θ is a real parameter and H is the Hopf invariant which we now define.

The conservation of J^μ licenses us to manufacture a "gauge potential" A^μ by the curl equation:

$$J^\mu = \epsilon^{\mu\nu\lambda} \partial_\nu A_\lambda \equiv \frac{1}{2} \epsilon^{\mu\nu\lambda} F_{\nu\lambda}. \quad (6)$$

A_μ is defined up to the gauge freedom $A_\mu \rightarrow A_\mu - \partial_\mu \Lambda$. Note that A_μ depends nonlocally on $n^\alpha(x)$. In the gauge $\partial A = 0$, we have $A_\mu = -\delta^{-2} \epsilon_{\mu\nu\lambda} \partial_\nu J_\lambda$. [An alternative construction is to write $A_\mu = iz^\dagger \partial_\mu z$. The U(1) phase rotation on z induces the gauge transformation on A_μ .] The Hopf invariant is defined by

$$H = -(1/4\pi) \int d^3x \epsilon^{\mu\nu\lambda} A_\mu F_{\nu\lambda} \\ = -(1/2\pi) \int d^2x A_\mu J^\mu. \quad (7)$$

H is obviously invariant under gauge transformation on A_μ . [We note that this is just the Abelian version of the topological term studied by Deser, Jackiw, and Templeton,⁸ but since H is gauge invariant, θ is not quantized. Furthermore, here H is to be regarded as a functional of $\vec{n}(\vec{x})$. In the language of Zumino, Wu, and Zee,¹² H is proportional to $\int \omega_3^0$. If $\hat{\mu}$ denotes a four-dimensional index then $\partial_\mu \epsilon^{\mu\nu\lambda} A_\mu F_{\nu\lambda} = \frac{1}{2} \epsilon^{\mu\nu\lambda} F_{\mu\lambda} F_{\nu\lambda}$, connecting the Hopf invariant to the chiral anomaly.]

Spatial rotation of a single skyrmion is equivalent to an isospin rotation and thus we evaluate H for the time-varying configuration $n_i \pm in_3 = e^{\pm i\alpha(t)} \times (\hat{n}_1 \pm i\hat{n}_2)(\vec{x})$, $n_3 = \hat{n}_3(\vec{x})$. {Strictly speaking, this defines a map of $S_2 \times [0, 1] \rightarrow S_2$.} It is not necessary to know the explicit form of n_a . From Eq. (3) we find

$$J_i = -(1/8\pi)(d\alpha/dt)\epsilon_{ijk} \partial_j n_3. \quad (8)$$

It suffices to know that $J_0(r) = \epsilon_{ijk} \partial_j A_i$ is a function of r to determine $A_i = -\epsilon_{ijk} x_k g(r)/r^2$ where

$g(r) = \int_0^r dr' r' J_0(r')$. This and Eq. (8) allow us to determine $A_0 = -(d\alpha/dt)n_3$. Inserting into Eq. (7) we find

$$H = g(\infty)n_3(\infty)[\alpha(T) - \alpha(0)]/2\pi = 1. \quad (9)$$

The skyrmion has angular momentum $\theta/2\pi$.

For a ferromagnet, θ should be determined by the microscopic theory underlying the phenomenological σ model.

It is easy to show that H is a homotopic invariant for $s_3 \rightarrow s_2$. Consider a map with $\vec{n}(\vec{x}, t = \infty)$ constant and a small deformation of \vec{n} leaving invariant $\vec{n}(\infty)$. Then

$$\begin{aligned} \delta J_\mu &= \epsilon_{\mu\nu\lambda} \partial_\nu \epsilon_{abc} 2n_a \delta n^b \partial_\lambda n^c \\ &= \epsilon_{\mu\nu\lambda} \partial_\nu \delta A_\lambda \end{aligned} \quad (10)$$

and we find

$$\delta H = (-1/2\pi) 2 \int d^3x \delta A_\mu J^\mu = 0.$$

There is a deep theorem¹³ which equates the Hopf invariant to the linking number between two curves in R^3 . To have a heuristic understanding of this, consider the maps $s_3 \rightarrow s_2$. The reverse image of a point in s_2 is a curve in s_3 which by a stereographic projection we can think of as a curve in R^3 (with ∞ identified as one point). Thus, for the basic map given explicitly above, $\vec{n} = (0, 0, 1)$ corresponds to the great circle $|z_1| = 1$, $z_2 = 0$ on s_3 , while $\vec{n} = (0, 0, -1)$ corresponds to $z_1 = 0$, $|z_2| = 1$. Write the real components of (z_1, z_2) as $(\cos\psi, \sin\psi \cos\theta, \sin\psi \sin\theta \cos\varphi, \sin\psi \sin\theta \sin\varphi)$ and stereographically project this point to $\vec{r}(\psi) \times (\cos\theta, \sin\theta \cos\varphi, \sin\theta \sin\varphi)$ in R^3 where $\vec{r}(\psi)$ ranges monotonically from ∞ to 0 as ψ ranges from 0 to π . We see that the curves corresponding to $\vec{n} = (0, 0, 1)$ and to $\vec{n} = (0, 0, -1)$ link once. The reader may find it amusing to work out the curves corresponding to other points.

Using this linking theorem, we can easily determine the spin and statistics of a skyrmion. Consider the following process in 2+1 dimensions. At some time create a pair of skyrmion and antiskyrmion and pull them apart. Rotate the skyrmion through 2π . Allow the pair to come together. Since at ∞ we have the physical vacuum this defines a map $s_3 \rightarrow s_2$. Were the skyrmion not rotated, the map would be homotopically trivial. Here, the corresponding map has Hopf invariant 1. The two curves traced out by two specific values of \vec{n} will be linked once as indicated in Fig. 1.

To determine the statistics obeyed by a skyrmion we consider a process in which we create two skyrmion-antiskyrmion pairs and subsequently

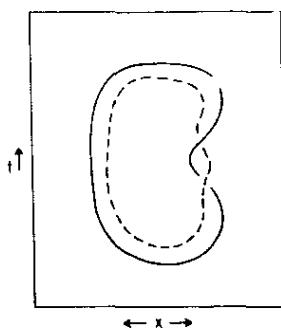


FIG. 1. The creation and annihilation of a skyrmion-antiskyrmion pair, with a 2π rotation of the skyrmion. The two curves correspond to $\vec{n} = (0, 0, 1)$ and $(1, 0, 0)$, say.

bring them to annihilation but after interchanging the two skyrmions. We see, by the maneuvering indicated in Fig. 2, that the linking number is 1 for this process. The map of $s_3 \rightarrow s_2$ corresponding to this process therefore has Hopf invariant 1. Thus, the skyrmion obeys exotic statistics which interpolates continuously between Bose and Fermi statistics as described in Ref. 6. (By the way, the alternative of directly computing the Hopf integral corresponding to rotating a pair of skyrmions through π appears to be quite difficult.) Note that the discussion there is for a gauge theory. Here, we do not have a gauge theory but, curiously, one can manufacture a gauge potential A_μ .

Given a map $f: s_k \rightarrow s_n$ one can always construct¹³ a map $\tilde{f}: s_{k+1} \rightarrow s_{n+1}$ (called the Freudenthal suspension of f) by $\tilde{f}(t, (1-t^2)^{1/2}x) = (t, (1-t^2)^{1/2}f(x))$ where $x \in s_k$ and $t \in [0, 1]$. This induces a homomorphism¹³ $F: \pi_k(s_n) \rightarrow \pi_{k+1}(s_{n+1})$ of the homotopy classes of f and \tilde{f} . Our discussion can thus be "suspended" into $(3+1)$ -dimensional space-time: $\pi_2(s_2) \rightarrow \pi_3(s_3) = z$ and $\pi_3(s_3) \rightarrow \pi_4(s_4) = z_2$. The first of these is an isomorphism, the second is onto: The suspension of a map $s^3 \rightarrow s^2$ to a map $s^4 \rightarrow s^3$ is nontrivial if and only if the map has odd Hopf invariant. Since $s_3 = \text{SU}(2)$ manifold, the homotopy $\pi_3(s_3)$ implies the existence of skyrmions in the $\text{SU}(2) \otimes \text{SU}(2)$ nonlinear σ model. The fact that $\pi_4(s_3) = z_2$ allows one to quantize the skyrmion as a spin- $\frac{1}{2}$ fermion as discussed by Witten.⁹ It is consistent with the standard three-space angular momentum analysis and with the well-known facts $\pi_1(\text{SO}(2)) = z$ and $\pi_1(\text{SO}(3)) = z_2$ that $\pi_4(s_3)$ is z_2 rather than z .

This material is based upon research supported in part by the National Science Foundation under Grant No. PHY77-27084, supplemented by funds from the National Aeronautics and Space Admini-

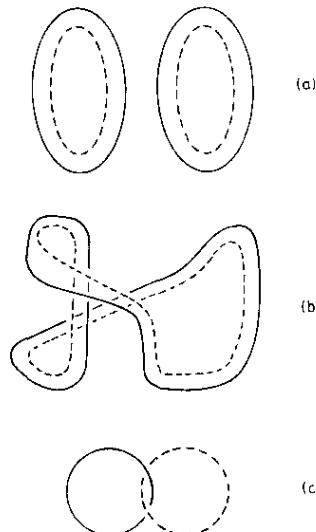


FIG. 2. (a) The creation and annihilation of two skyrmion-antiskyrmion pairs. (b) The process in (a) but with an interchange of the two skyrmions. (c) The two curves in (b) after a homotopic deformation.

stration.

^(a)On leave of absence 1983-1984.

¹A. A. Belavin and A. M. Polyakov, Pis'ma Zh. Eksp. Teor. Fiz. 22, 503 (1975) [JETP Lett. 22, 245 (1975)].

²T. H. R. Skyrme, Proc. Roy. Soc. London, Ser. A 247, 260 (1958).

³For a review of the material in this introductory paragraph, see R. Rajaraman, *Solitons and Instantons* (North-Holland, Amsterdam, 1982).

⁴Some aspects of this subject appear to have been anticipated in the remarkable paper of D. Finkelstein and J. Rubinstein, J. Math. Phys. 9, 1762 (1968).

⁵P. Hasenfratz, Phys. Lett. 85B, 338 (1979); J. Schonfeld, Nucl. Phys. B185, 157 (1981).

⁶F. Wilczek, Phys. Rev. Lett. 48, 1144 (1982), and 49, 957 (1982).

⁷G. 't Hooft, Phys. Rev. Lett. 37, 8 (1976); C. Callan, R. Dashen, and D. Gross, Phys. Lett. 63B, 334 (1976); R. Jackiw and C. Rebbi, Phys. Rev. Lett. 37, 172 (1976).

⁸J. Schonfeld, Ref. 5; S. Deser, R. Jackiw, and S. Templeton, Phys. Rev. Lett. 48, 975 (1982), and Ann. Phys. (N.Y.) 140, 372 (1982); see also Y.-S. Wu and A. Zee, to be published.

⁹E. Witten, Nucl. Phys. B223, 422, 433 (1983).

¹⁰A. P. Balashandran, V. P. Nair, and C. G. Trahern, Phys. Rev. Lett. 49, 1124 (1982).

¹¹J. Wess and B. Zumino, Phys. Lett. 37B, 95 (1971).

¹²B. Zumino, Y.-S. Wu, and A. Zee, to be published.

¹³For example, P. J. Hilton, *An Introduction to Homotopy Theory* (Cambridge Univ. Press, Cambridge, England, 1953), Chap. VI.

Multiparticle Quantum Mechanics Obeying Fractional Statistics

Yong-Shi Wu

Department of Physics, University of Washington, Seattle, Washington 98195

(Received 2 April 1984)

We obtain the rule governing many-body wave functions for particles obeying fractional statistics in two (space) dimensions. It generalizes and continuously interpolates the usual symmetrization and antisymmetrization. Quantum mechanics of more than two particles is discussed and some new features are found.

PACS numbers 03.65.Ca, 03.65.Ge, 05.30.-d

In two (space) dimensions, there are allowed to be particles of fractional angular momentum or spin.^{1,2} If there is a generalized spin-statistics connection, such particles are expected to have unusual (fractional) statistics which continuously interpolates between the normal bosons and fermions. (An example for such interpolation is known in one dimension.³) The intriguing problem of how it works is interesting both from the viewpoint of theoretical principles and from the prospect of physical applications. A possible relevance of fractional statistics to the quantized Hall effect has been recently suggested.⁴

Two simple models have been proposed for particles obeying fractional statistics by Wilczek^{1,5} and Wilczek and Zee.² (See also Ref. 6.) Two-particle quantum mechanics was analyzed in detail. A low-density expansion of the partition function interpolating the standard statistics was obtained. As pointed out in these papers, Feynman's path-integral formulation is a good starting point. However, the formalism in terms of wave functions may

be practically more convenient. An immediate problem is the general rule governing the many-body wave functions, namely how to generalize the usual rule to obtain a continuous interpolation between symmetrization and antisymmetrization. In this note I answer this question by deriving the desired rule in the two models mentioned above. As an application, I discuss the quantum mechanics of three particles, not yet touched in the literature. Some new features are found which are not present in the two-particle case.

Anyons revisited.—Following Wilczek,⁵ I denote composites formed from charged particles and magnetic flux tubes as anyons, since their spin

$$\Delta = q\Phi/2\pi = \theta/2\pi \quad (1)$$

can take any real values. Here $-q$ is the charge and Φ the flux. That⁵ interchange of two anyons leads to a phase $e^{i\theta}$ is an indication of the fractional statistics. We here consider quantum mechanics for more than two anyons.

The Hamiltonian for a charged particle orbiting around a flux tube can be written as

$$H_0 = \frac{1}{2m_q} \left[-i \frac{\partial}{\partial \vec{r}_q} + q \vec{A}(\vec{r}_q - \vec{r}_f) \right]^2 + \frac{1}{2m_f} \left[-i \frac{\partial}{\partial \vec{r}_f} - q \vec{A}(\vec{r}_q - \vec{r}_f) \right]^2. \quad (2)$$

Here we consider the limit in which the size of the flux tube can be neglected. \vec{r}_q and \vec{r}_f are two-dimensional vectors. Let us assume that the flux tube has a finite effective mass m_f in two dimensions. The form (2) has the advantage that the effect of the interaction is confined to the wave function in the relative coordinate. In a regular gauge the vector potential is

$$q \vec{A}(\vec{r}_q - \vec{r}_f) = -q \vec{A}(\vec{r}_f - \vec{r}_q) = (\theta/2\pi)[\vec{n} \times (\vec{r}_q - \vec{r}_f)]/|\vec{r}_q - \vec{r}_f|^2 \quad (3)$$

(with \vec{n} being the unit vector normal to the two-dimensional plane), and the wave function is single-valued everywhere.

Now we proceed to consider n identical anyons and neglect the electrostatic forces between them (i.e., consider the limit $q \rightarrow 0$ with $\theta = q\Phi$ fixed). The charged particle in each anyon feels the vector potential of the flux tube in the other. Using the Hamiltonian (1) and applying a procedure similar to that in Goldhaber⁷ for the charge-monopole composites, one finds that the anyon-anyon potential is equivalent to that of a charge interacting with twice the flux in one flux tube; namely

$$H = \sum_{i=1}^n \frac{1}{2m_a} \left[-i \frac{\partial}{\partial \vec{r}_i} + 2q \sum_{j \neq i} \vec{A}(\vec{r}_i - \vec{r}_j) \right]^2. \quad (4)$$

Let us adopt Eq. (3) for $\vec{A}(\vec{r}_i - \vec{r}_j)$ in the regular gauge in which the wave function ψ is single-valued as in the one anyon case. To eliminate the long-range vector potential between anyons, we make the gauge transformation

$$\psi'(\vec{r}_1, \dots, \vec{r}_n) = \prod_{i < j} \exp\left(i\frac{\theta}{\pi}\phi_{ij}\right) \psi(\vec{r}_1, \dots, \vec{r}_n), \quad (5)$$

where ϕ_{ij} is the azimuthal angle of the relative vector $\vec{r}_i - \vec{r}_j$. Now the new wave function ψ' satisfies the free Schrödinger equation with no vector potential.

At first sight the multivaluedness of the new wave function ψ' seems to be very discomforting. One can manage to avoid it by imposing appropriate boundary conditions for ψ' on certain cuts in the two-dimensional plane⁵ or formulating quantum mechanics on sections on fiber bundles.⁸ However, these two methods are very hard to put into practice for more than two anyons. Actually, nothing is wrong with the multivaluedness of the wave function (5). The modulus squared, $|\psi'|^2$, is single-valued, and the multivalued phase factors are just right to keep track of the Aharonov-Bohm effect.⁹ In my opinion once one understands the need for extending the notion of a wave function (i.e., not requiring it to be necessarily 2π periodic in ϕ_{ij}), there is no difficulty in accepting and directly using the multivalued wave function (5) as everybody does with the double-valued spinors in three dimensions.¹⁰

By use of the complex coordinates $z_i = x_i + iy_i$ and $z_i^* = x_i - iy_i$, instead of $\vec{r}_i = (x_i, y_i)$, the wave function (5) can be put into a more elegant form¹¹:

$$\psi'(z_i, z_i^*) = \prod_{i < j} (z_i - z_j)^{\theta/\pi} f(z_i, z_i^*), \quad (6)$$

with $f(z_i, z_i^*) = (r_{ij})^{-\theta/\pi} \psi(z_i, z_i^*)$ single-valued. f is totally symmetric (antisymmetric) in the pairs (z_i, z_i^*) , if all the fields describing the flux tube and charged particle are bosonic (if the charged particle is fermionic). The equation (6) is the desired rule for many-body wave functions obeying θ statistics.

Solitons in point approximation.—The solitons in the $(2+1)$ -dimensional O(3) nonlinear sigma model, with a topological action, also provide a model for particles with fractional spin and statistics.^{2,6} When widely separated solitons are approximately treated as point particles, the topological term (with the parameter θ) leads to an additional term

$$S' = \int dt L', \quad L' = (-\theta/\pi)(d/dt) \sum_{i < j} \phi_{ij}, \quad (7)$$

to the ordinary action $S_0 = \int dt \frac{1}{2} m \sum_i \dot{\vec{r}}_i^2$. While this term does not affect the equation of motion, it determines the statistics of the particles via path integral.

When one goes from path integral to wave functions, the term (7) also leads to the rule (6) for many-body wave functions associated with usual Hamiltonian containing no peculiar interactions. In fact, the change of ϕ_{ij} can be always written as

$$\phi_{ij}(t)|_{t'}^{t''} = 2\pi n_{ij} + \phi_{ij}' - \phi_{ij}, \quad (8)$$

with $0 \leq \phi_{ij}' - \phi_{ij} < 2\pi$. Thus, the propagator in the n -particle configuration space is a sum of "partial amplitudes," each corresponding to a distinct class of paths having the same winding numbers $\{n_{ij}\}$:

$$K(\vec{r}_i'', t''; \vec{r}_i', t') = \exp\left[-i\frac{\theta}{\pi} \sum_{i < j} (\phi_{ij}'' - \phi_{ij}')\right] \sum_{n_{ij}} \exp(-i2\theta \sum_{i < j} n_{ij}) \int_{\vec{r}_i'}^{\vec{r}_i''} [\mathcal{D}\vec{r}_i(t)]_{n_{ij}} \exp(iS_0). \quad (9)$$

As usual, a single-valued wave function $\psi(\vec{r}_i, t)$ can be introduced such that

$$\psi(\vec{r}_i'', t'') = \int d\vec{r}_i' K(\vec{r}_i'', t''; \vec{r}_i', t') \psi(\vec{r}_i', t'). \quad (10)$$

We can eliminate the sum in Eq. (9) by introducing a new wave function

$$\tilde{\psi}(\vec{r}_i, t) = \exp\left(i\frac{\theta}{\pi} \sum_{i < j} \phi_{ij}\right) \psi(\vec{r}_i, t). \quad (11)$$

Then, corresponding to Eq. (10), now we have

$$\begin{aligned} \tilde{\psi}(\vec{r}_i'', t'') \\ = \int d\vec{r}'_i \tilde{K}(\vec{r}_i'', t''; \vec{r}'_i, t') \tilde{\psi}(\vec{r}'_i, t'), \end{aligned} \quad (12)$$

$$\begin{aligned} \tilde{K}(\vec{r}_i'', t''; \vec{r}'_i, t') \\ = \int_{\vec{r}'_i}^{\vec{r}_i''} [\mathcal{D}\vec{r}_i(t)] \exp(iS_0). \end{aligned} \quad (13)$$

$$H = \frac{m}{2} \sum_i \dot{\vec{r}}_i^2, m \dot{\vec{r}}_i = \vec{p}_i - \frac{\theta}{\pi} \sum_{i < j} \frac{(\vec{r}_i - \vec{r}_j) \times \vec{n}}{|\vec{r}_i - \vec{r}_j|^2}.$$

Here \vec{p}_i is the canonical momentum conjugate to \vec{r}_i . It is easy to see that H is the same as given by Eq. (4) together with Eq. (3). We can repeat the same procedure in the last section to arrive at Eq. (6). However, the argument given from Eq. (8) to Eq. (14) has the advantage that it elucidates the relationship between our wave functions and the path integral formulation.

Properties of the wave function (6).—Equation (5) or the rule (6) is invariant under $\theta \rightarrow \theta + 2\pi$; i.e., fractional statistics is 2π periodic in θ , in agreement with the well-known periodicity of the Aharonov-Bohm effect in the flux or that of the θ parameter in the topological action.

When $\theta = 0$ and π , the rule (6) coincides with the standard symmetric or antisymmetric rule. For intermediate θ it gives a continuous interpolation between the two extreme cases. However, when $\theta \neq 0, \pi$, the many-body wave functions are not of the form of products of single-particle wave functions. So generally we expect that the physical quantities of a system of many particles are not simply related to those for one particle.

When $n = 2$, from Eq. (6) it is easy to recover the condition^{1,5}

$$\psi'(\phi_{12} \pm \pi) = e^{\pm i\theta} \psi'(\phi_{12}). \quad (15)$$

For $n \geq 3$, Eq. (6) exhibits complicated behavior

$$H\psi = E\psi, \quad H = -\frac{1}{2} \sum_{i=1}^n \frac{\partial^2}{\partial z_i \partial z_i^*} + \frac{1}{2} \omega^2 \sum_{i=1}^n z_i z_i^*,$$

where ψ satisfies the rule (6) with f totally symmetric. (We have omitted the prime on ψ).

The $n = 2$ case has been analyzed in Refs. 5 and 12. In our approach we recover the complete set of solutions as follows:

$$\psi = W^{|L|} w^{l+2\Delta} L_N^{(|L|)} (2\omega Z Z^*) L_n^{(|l+2\Delta|)} (\frac{1}{2}\omega z z^*) \exp[\frac{1}{2}\omega(z_1 z_1^* + z_2 z_2^*)]. \quad (17)$$

$$E = (2N + |L| + 2n + |l+2\Delta| + 2)\omega. \quad (18)$$

where $N, n \geq 0$ are principal quantum numbers for the center-of-mass and relative oscillators respectively; L , $l+2\Delta$ are angular momenta in the center-of-mass and relative coordinates. (l must be even.) $L_M^{(m)}(x)$ are

Note that the wave function $\tilde{\psi}(\vec{r}_i, t)$ is single-valued on the universal covering space (or Riemann surface) of the n -particle configuration space. The integration over \vec{r}'_i in Eq. (12) is taken on this covering space. By use of the complex coordinates, it is easy to recover Eq. (6) from Eq. (11).

Another way to derive the same result is the following. The Hamiltonian corresponding to $L_0 + L'$ is

$$(14)$$

under permutation or interchange of the positions of particles. Complication occurs even when we exchange only two particles in the presence of a third particle. We have to specify along what loop particle 1 moves from \vec{r}_1 to \vec{r}_2 and particle 2 from \vec{r}_2 to \vec{r}_1 . The resulting phase change will depend on whether the "spectator" 3 is enclosed inside this loop or not. This situation is a reflection of the fact that the configuration space of identical particles is multiply connected. It is the origin of the difficulties pointed out in Refs. 5 and 6 in dealing with more than two particles. The acceptance and direct use of the multivalued wave functions (6) make the many-particle problem accessible to approach, since the complications mentioned above have been simply built into the factors $\prod_{i < j} (z_i - z_j)^{\theta/\pi}$.

Physically, the long-range interactions due to θ statistics are coded in the factors $\prod_{i < j} (z_i - z_j)^{\theta/\pi}$. Moreover, these factors imply the existence of angular momentum barriers between any pair of particles when $\theta \neq 0$. Thus the many-body wave functions are expected to vanish when any two of the particles coincide (if $\theta \neq 0$), although the particles are not fermions for $\theta \neq \pi$.

Three particles, harmonic well.—As an application let us use the many-body wave functions (6) to attack the problem of three identical particles in a harmonic potential. The Schrödinger equation (for n particles with $m = 1$) is

$$(16)$$

$$113$$

the Laguerre polynomials.¹³ We have used the following notation for brevity: $Z = \frac{1}{2}(z_1 + z_2)$, $z = z_1 - z_2$ and

$$W = \begin{cases} Z & \text{if } L > 0, \\ Z^* & \text{if } L < 0, \end{cases} \quad w = \begin{cases} z & \text{if } l + 2\Delta > 0, \\ z^* & \text{if } l + 2\Delta < 0. \end{cases} \quad (19)$$

Since θ appears only in the form of $|l + 2\Delta|$, the 2π periodicity of θ is made clear. It is also easy to verify the continuous interpolation between the spectrum (including degeneracies) of bosons and that of fermions when θ varies from 0 to π .^{12,14}

For $n = 3$, we have obtained the following solutions for $0 \leq \theta < \pi$:

$$\psi = [(z_1 - z_2)(z_1 - z_3)(z_2 - z_3)]^{\theta/\pi} \exp\{-\frac{1}{2}\omega r^2\} P, \quad (20)$$

$$P = (z_1 + z_2 + z_3)^L (z_1 - z_2)^l (2z_1 - z_2 - z_3)^m L_{N_1}^{(L)} (\frac{1}{3}\omega R^2) L_{N_2}^{(l+3m+6\Delta-5)} (\frac{1}{3}\omega \rho^2) + \text{symmetrization}, \quad (20')$$

$$E = (2N_1 + 2N_2 + L + l + m + 6\Delta + 3)\omega, \quad (20'')$$

where all N_1, N_2, L, m, l are nonnegative integers, and l, m such that after symmetrization P does not become identically zero. Moreover, $R^2 = |z_1 + z_2 + z_3|^2$, $r^2 = \sum |z_i|^2$,

$$\rho^2 = |2z_1 - z_2 - z_3|^2 + \text{cyclic permutation}.$$

We note that the parity transformation $z_i \rightarrow z_i^*$ and $\theta \rightarrow -\theta$ is a good symmetry of the equation (16) and the rule (6). So applying it on the solutions (20) will lead to more solutions (with l, m such that ψ has no singularities at $z_i^* = z_j^*$). We know that this set of solutions does not exhaust those of the problem; e.g., the three-fermion ground state is missing when $\theta = \pi$.

Even so, we are able to see some important features not present in the solutions of two particles. First, for sufficiently small θ , the ground-state energy is $E_0 = (3 + 3\theta/\pi)\omega$. For n particles, it is $E_0 = [n + n(n-1)\theta/2\pi]\omega$. Thus, the n dependence of E_0 has a quadratic part which looks like two-body interaction energy. Second, when $\theta = \pi$ the above energy level moves to 6ω , which exceeds the energy of the three-fermion ground state $E'_0 = 5\omega$. So when θ varies continuously from 0 to π , there must be level crossing and, therefore, the emergence of new ground states at certain values of θ . This effect may lead to interesting phenomena in realistic systems obeying θ statistics when θ can vary under certain circumstances.

To conclude, I stress that though the rule (6) is derived in two concrete models, it is generally true for any fractional statistics in two dimensions, whatever its origins. This will be confirmed in a model-independent formulation in a forthcoming paper.¹⁵

The author thanks M. Baker, D. Boulware, L. Brown, R. Tao, F. Wilczek, and A. Zee for useful discussions and comments. This work was supported in part by the U. S. Department of Energy.

¹F. Wilczek, Phys. Rev. Lett. **48**, 1144 (1982). See also M. Peshkin, Phys. Rep. **80**, 376 (1982).

²F. Wilczek and A. Zee, Phys. Rev. Lett. **51**, 2250 (1983).

³C. N. Yang and C. P. Yang, J. Math. Phys. **10**, 1115 (1969). They have shown how a Bose gas with repulsive two-body delta-function interactions becomes a free Fermi gas for infinite coupling constant.

⁴B. Halperin, Phys. Rev. Lett. **52**, 1583, 2390(E) (1984).

⁵F. Wilczek, Phys. Rev. Lett. **49**, 957 (1982).

⁶F. Wilczek and A. Zee, Institute of Theoretical Physics, University of California, Santa Barbara Report No. NSF-ITP-84-25, 1984 (to be published).

⁷A. S. Goldhaber, Phys. Rev. Lett. **36**, 1122 (1976), and **49**, 905 (1982).

⁸T. T. Wu and C. N. Yang, Phys. Rev. D **12**, 3845 (1975).

⁹Y. Aharonov and D. Bohm, Phys. Rev. **115**, 485 (1959).

¹⁰We note here that infinite multivaluedness of a wave function can happen only in two dimensions, since $\pi_1(SO(2)) = Z$. In three-space only double-valuedness is allowed because $\pi_1(SO(3)) = Z_2$.

¹¹Here by the notation $\psi(z, z_i^*)$, the set $\{(z, z_i^*)\}$, $i = 1, \dots, n\}$ is understood. Special wave functions of this form have appeared in Ref. 4. Here we proved that Eq. (6) is the general form of many-body functions obeying fractional statistics.

¹²J. Leinaas and J. Myrheim, Nuovo Cimento Soc. Ital. Fis. **B37**, 1 (1977).

¹³See, e.g., *Encyclopedic Dictionary of Mathematics*, edited by S. Iyanaga and Y. Kawada (MIT Press, Cambridge, Mass., 1977), Appendix A, Table 20 VI.

¹⁴The continuous interpolation between the two-body scattering amplitudes of bosons and those of fermions is being discussed by F. Wilczek and A. Zee (private communication).

¹⁵Y. S. Wu, to be published.

STATISTICAL MECHANICS OF ANYONS

Daniel P. AROVAS

Department of Physics, University of California, Santa Barbara, CA 93106, USA

Robert SCHRIEFFER and Frank WILCZEK

*Department of Physics, University of California
and*

Institute for Theoretical Physics, Santa Barbara, CA 93106, USA

A. ZEE*

Institute for Advanced Study, Princeton, NJ 08540, USA

Received 31 July 1984

We study the statistical mechanics of a two-dimensional gas of free anyons – particles which interpolate between Bose-Einstein and Fermi-Dirac character. Thermodynamic quantities are discussed in the low-density regime. In particular, the second virial coefficient is evaluated by two different methods and is found to exhibit a simple, periodic, but nonanalytic behavior as a function of the statistics determining parameter.

In two space dimensions, a continuous family of quantum statistics interpolating between bosons and fermions is possible [1, 2]. Two examples of particles obeying exotic statistics have been discussed. The soliton of the $(2+1)$ -dimensional O(3) nonlinear σ -model has a spin which is neither integral nor half-odd integral, and obeys a statistics which is neither Bose-Einstein nor Fermi-Dirac [3]. In condensed matter, one finds that the Laughlin quasiparticles [4] in the anomalous quantum Hall effect system also possess fractional charge and fractional statistics [5], a result recently derived from the adiabatic theorem [6].

We first discuss a method [1] by which the statistics of a two-dimensional system of charged particles can be changed (continuously) via the introduction of a fictitious “statistical gauge field.” It is well known that the wave function of a charged particle interacting with a magnetic flux tube will acquire a phase change due to the motion of the particle. If the charge is $e^* = \nu e$ and the flux is

* On leave from the University of Washington 1983–1984.

$\phi = \alpha\phi_0 = \alpha hc/e$, a complete revolution will induce a phase change of $e^{i\Delta\gamma}$ with $\Delta\gamma = 2\pi\alpha\nu$. The basic idea of this method is then to introduce a flux tube at the position of each particle whose magnitude ϕ will determine the phase associated with relative particle motion. The gauge field associated with this flux is nondynamical, i.e., its evolution is completely determined by the motion of the particles. It is expected that physical quantities are periodic in the statistics determining parameter α .

In the symmetric gauge, the vector potential due to a flux tube of magnitude $\phi = \alpha\phi_0$ takes the form

$$\mathbf{A}(\mathbf{r}) = \frac{\alpha\phi_0}{2\pi} \hat{z}x \frac{\mathbf{r}}{r^2} = \frac{\alpha\phi_0}{2\pi} \frac{\hat{\theta}}{r} = \frac{\alpha\phi_0}{2\pi} \nabla\theta. \quad (1)$$

The many-particle generalization we seek is therefore*

$$\mathbf{A}_\phi(\mathbf{r}_i) = \frac{\alpha\phi_0}{2\pi} \hat{z}x \sum_j' \frac{\mathbf{r}_{ij}}{r_{ij}^2} = \frac{\alpha\phi_0}{2\pi} \sum_j' \nabla_i \theta_{ij}, \quad (2)$$

where $\theta_{ij} = \tan^{-1}((y_j - y_i)/(x_j - x_i))$ is the relative angle between i and j , and the prime on the sum indicates that the term $j = i$ is to be excluded. This leads to the following many-body hamiltonian:

$$H(\alpha) = \sum_i \frac{1}{2m} \left(\mathbf{p}_i - \frac{e}{c} \mathbf{A}_0(\mathbf{r}_i) - \frac{e}{c} \mathbf{A}_\phi(\mathbf{r}_i) \right)^2 + V(\mathbf{r}_1, \dots, \mathbf{r}_N), \quad (3)$$

where \mathbf{A}_0 is the physical vector potential, if present. Suppose one knows an eigenfunction ψ_0 of the bare hamiltonian $H_0 \equiv H(0)$ with energy E_0 . Then since

$$\left[\mathbf{p}_i, \exp\left(i \frac{\alpha\phi_0}{2\pi} \sum_j' \theta_{ij}\right) \right] = \frac{\alpha\phi_0}{2\pi} \sum_j' \nabla_i \theta_{ij} = \mathbf{A}_\phi(\mathbf{r}_i), \quad (4)$$

we see that $\psi_\alpha \equiv \exp(i\alpha\sum_{\text{pairs}} \theta_{ij})\psi_0$ is an eigenfunction of $H(\alpha)$ also with energy E_0 . The problem is that the function ψ_α will not in general be a single-valued function of its arguments.

For the two-particle problem with no external potential and no interparticle interactions (free particles), the situation becomes eminently tractable. Recall that the wave function can be decomposed into a product $\psi_0(\mathbf{r}, \mathbf{R}) = \chi(\mathbf{R})\xi(\mathbf{r})$ where \mathbf{R} is the center-of-mass position and \mathbf{r} is the relative coordinate vector. We find $\chi(\mathbf{R}) = e^{i\mathbf{K} \cdot \mathbf{R}}$, $\xi(\mathbf{r}) = e^{im\theta} J_{|m|}(kr)$, $E = \hbar^2 K^2 / 4M + \hbar^2 k^2 / M$ (M = particle mass). Imposition of Bose (Fermi) statistics then requires that m be even (odd). The introduction of a statistical gauge field then provides us with new eigenfunctions $\psi_\alpha(\mathbf{r}, \mathbf{R}) = e^{i\alpha\theta} \psi_0(\mathbf{r}, \mathbf{R}) = e^{i\mathbf{K} \cdot \mathbf{R}} e^{i(m+\alpha)\theta} J_{|m|}(kr)$. Here, Bose (Fermi) statistics re-

* Since there is flux-charge interaction as well as charge-flux interaction, $\alpha = 2\phi/\phi_0$ is now twice the flux per tube in units of the Dirac quantum.

quires that $m + \alpha \equiv l$ be even (odd) and we obtain wave functions

$$\psi_\alpha(\mathbf{r}, \mathbf{R}) = e^{i\mathbf{K} \cdot \mathbf{R}} e^{il\theta} J_{|l-\alpha|}(kr). \quad (5)$$

If we now introduce a circular boundary at some radius \tilde{R} , we find that the allowed energies are

$$\varepsilon_{l,n} = \hbar^2 x_{|l-\alpha|,n}^2 / M\tilde{R}^2, \quad J_\nu(x_{\nu,n}) = 0. \quad (6)$$

Hence, choosing α to be an odd integer merely shifts the energy spectrum from Bose-like to Fermi-like. Note also that the spectrum is periodic in α with period $\Delta\alpha = 2$. This is in fact true for the N -particle system, although explicit (single-valued) wave functions are difficult to obtain due to the fact that there are $\frac{1}{2}N(N-1)$ relative angles and only $(N-1)$ non-CM angular degrees of freedom. For $N=2$, these numbers are identical.

This result, eq. (6), can be used to evaluate the second virial coefficient. In two dimensions, it is easily shown that [7]*

$$B(T) = \frac{1}{2}A - A^{-1}\lambda_T^4 Z_2, \quad (7)$$

where A is the area of the system $\lambda_T = (2\pi\hbar^2/MkT)^{1/2}$ is the thermal wavelength, and $Z_2 = \text{Tr}e^{-\beta H_2}$ is the two-particle partition function. The virial expansion is an expansion of the equation of state in the density n : $P = nkT[1 + Bn + Cn^2 + \dots]$. In performing the trace to obtain Z_2 , the center-of-mass freedom is trivially separated, yielding a factor $Z_2 = 2A\lambda_T^{-2}\tilde{Z}_2$, where \tilde{Z}_2 is now the single particle partition function for the relative coordinate problem: $\tilde{Z}_2 = \text{Tr}_{\text{rel}}e^{-\beta H_{\text{rel}}}$. This will again be area divergent, and it is therefore convenient to calculate the virial coefficient $B(\alpha, T)$ in terms of a known quantity, i.e., $B(2j, T) = -\frac{1}{4}\lambda_T^2$ or $B(2j+1, T) = +\frac{1}{4}\lambda_T^2$, the familiar result for Bose and Fermi systems, respectively ($j \in \mathbb{Z}$). Thus, we obtain

$$B(\alpha', T) - B(\alpha, T) = 2\lambda_T^2 [\tilde{Z}_2(\alpha) - \tilde{Z}_2(\alpha')]. \quad (8)$$

We now appeal to the result (6). Clearly $B(\alpha, T)$ must be periodic in α with period $\Delta\alpha = 2$. We will take our original particles to have Bose statistics and expand about even and odd values of α . For $\alpha = 2j + \delta$, $|\delta| < 1$, corresponding to quasi-bosons, the allowed values of $|l - \alpha|$ are $|\delta|, 2 \pm \delta, 4 \pm \delta$, etc. For $\alpha = 2j + 1 + \delta$, $|\delta| < 1$, corresponding to quasi-fermions, the allowed values of $|l - \alpha|$ are $1 \pm \delta, 3 \pm \delta$, etc. Since B must be independent of the cutoff \tilde{R} in the limit $\tilde{R} \rightarrow \infty$, and since \tilde{R}

* To be precise, we should write $B(T) = A[\frac{1}{2} - Z_2/Z_1^2]$ with $Z_1 = \text{Tr}e^{-\beta H_1}$, the single particle partition function. With no external fields, we have $H_1 = p^2/2M$ and $Z_1 = A\lambda_T^{-2}$, which then yields eq. (7)

appears only in the combination $MkT\tilde{R}^2/\hbar^2$, it is desirable to rescale $\tilde{R} \rightarrow \sqrt{\hbar^2/MkT}\tilde{R}$. Expanding about the Fermi point, we find that (8) and (6) give

$$B(2j+1+\delta, T) = \frac{1}{4}\lambda_T^2 + 2\lambda_T^2 \times \lim_{\tilde{R} \rightarrow \infty} \sum_{l=1}^{\infty} \sum_{n=1}^{\infty} \text{odd} [2e^{-(x_{l,n}/\tilde{R})^2} - e^{-(x_{l+\delta,n}/\tilde{R})^2} - e^{-(x_{l-\delta,n}/\tilde{R})^2}]. \quad (9)$$

The factor in brackets resembles a second derivative. By expanding in δ , one can then perform the l -sum by means of the celebrated Euler-MacLaurin formula [8]. This leaves

$$B(2j+1+\delta, T) = \frac{1}{4}\lambda_T^2 - 2\lambda_T^2\delta^2 \lim_{\tilde{R} \rightarrow \infty} \tilde{R}^{-2} \sum_{n=1}^{\infty} x_{1,n} \left. \frac{\partial x_{1,s,n}}{\partial s} \right|_{s=0} e^{-(x_{1,n}/\tilde{R})^2}. \quad (10)$$

As $n \rightarrow \infty$, $x_{1,n} \rightarrow \infty$ and

$$\frac{\partial x_{\nu,n}}{\partial \nu} = \frac{1}{J_{\nu+1}(x_{\nu,n})} \frac{\partial J_{\nu}(x_{\nu,n})}{\partial \nu} \rightarrow \frac{1}{2}\pi.$$

The value of n at which this approximation $x_{\nu,n} \sim \frac{1}{2}\nu\pi + n\pi - \frac{1}{4}\pi$ becomes valid is n_0 , say, which is certainly independent of \tilde{R} . The sum will then be completely dominated by the terms $n_0 \leq n < \infty$, the beginning terms being suppressed by the \tilde{R}^{-2} factor. Making this replacement, and writing $\sum_n \rightarrow \int dx_{1,n}/\pi$, we obtain

$$B(2j+1+\delta, T) = \frac{1}{4}\lambda_T^2 - \lambda_T^2\delta^2 \lim_{\tilde{R} \rightarrow \infty} \tilde{R}^{-2} \int_{n_0\pi}^{\infty} dx x e^{-(x/\tilde{R})^2} = \frac{1}{4}\lambda_T^2 - \frac{1}{2}\delta^2\lambda_T^2. \quad (11)$$

One can check that all other terms in the expansion in δ and in other approximations employed are formally of order \tilde{R}^{-1} as $\tilde{R} \rightarrow \infty$. Thus, we predict

$$B(2j+1+\delta, T) = \frac{1}{4}\lambda_T^2(1-2\delta^2)_{\text{per}}, \quad (12)$$

where the subscript indicates that we are to extend this function for $|\delta| > 1$ in a periodic fashion. The complete result has a cusp at Bose values $\alpha = 2j$ due to the required periodic extension. This in fact follows from the general formula (8). The only difference between the Fermi and Bose expansions is the existence of the

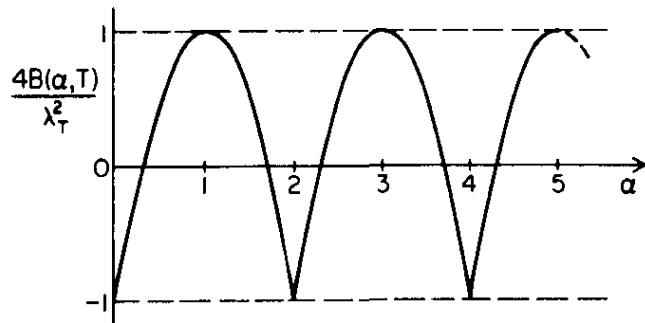


Fig. 1. The second virial coefficient $B(\alpha, T)$ as a function of the statistics determining parameter α (T fixed).

$|l - \alpha| = |\delta|$ term:

$$\begin{aligned} B(2j + \delta, T) &= -\frac{1}{4}\lambda_T^2 - \frac{1}{2}\delta^2\lambda_T^2 + 2\lambda_T^2 \lim_{\tilde{R} \rightarrow \infty} \sum_{n=1}^{\infty} [e^{-(x_{0,n}/\tilde{R})^2} - e^{-(x_{|\delta|,n}/\tilde{R})^2}] \\ &= -\frac{1}{4}\lambda_T^2 + |\delta|\lambda_T^2 - \frac{1}{2}\delta^2\lambda_T^2. \end{aligned} \quad (13)$$

This is exactly the form predicted above. Thus, we obtain a picture of $B(\alpha, T)$ for fixed T as in fig. 1.

This result is also derivable from a path integral approach. The general lagrangian for the many-body system is

$$L = \sum_i \frac{1}{2}M\dot{r}_i^2 + \alpha \sum_{\text{pairs}} \dot{\theta}_{ij}. \quad (14)$$

For a system of bosons, the partition function takes the form

$$Z_N = \frac{1}{N!} \int d^2r_1 \dots d^2r_N \sum_P \langle r_1 \dots r_N | e^{-\beta H_N} | P r_1 \dots P r_N \rangle, \quad (15)$$

which may be cast into a path integral form, as done by Feynman [9]. Again, the case $N = 2$ is considered, and the CM contribution is directly integrated out. This leaves

$$\begin{aligned} \tilde{Z}_2 &= \frac{1}{2} \int d^2r [\langle r | e^{-\beta H_{\text{rel}}} | r \rangle + \langle r | e^{-\beta H_{\text{rel}}} | -r \rangle], \\ L_{\text{rel}} &= \frac{1}{4}M\dot{r}^2 + \alpha\dot{\theta}. \end{aligned} \quad (16)$$

The $\dot{\theta}$ term in L introduces a winding number-dependent phase in the path integral. This problem is in fact equivalent to the Bohm-Aharonov effect [10], the path integral formulation of which was studied extensively by Gerry and Singh [11–14].

The matrix element

$$K(\mathbf{r}, \mathbf{r}'; \tau) = \int D\mathbf{r}(t) e^{i\hbar^{-1} \int_0^\tau dt' L(t')} = \langle \mathbf{r} | e^{-iH\tau/\hbar} | \mathbf{r}' \rangle, \\ \mathbf{r}(t=0) = \mathbf{r}, \quad \mathbf{r}(t+\tau) = \mathbf{r}', \quad \beta = i\tau, \quad (17)$$

can be decomposed into a sum over contributions of different homotopy sectors, with $\theta' - \theta \equiv \phi + 2\pi n$:

$$K(\mathbf{r}, \mathbf{r}'; \tau) = \sum_{n=-\infty}^{\infty} e^{2\pi i n \alpha} \bar{K}_n(\mathbf{r}, \mathbf{r}'; \tau), \quad (18)$$

$$\bar{K}_n(\mathbf{r}, \mathbf{r}'; \tau) = \int D\mathbf{r}(t) e^{i\hbar^{-1} \int_0^\tau dt' L(t')} \delta(\theta' - \theta - \phi - 2\pi n), \quad (19)$$

$$\begin{aligned} \bar{K}_n(\mathbf{r}, \mathbf{r}'; \tau) &= \frac{M}{4\pi\hbar i\tau} \exp\left[\frac{-M}{4\hbar i\tau}(r^2 + r'^2)\right] \\ &\times \int_{-\infty}^{\infty} d\lambda e^{i\lambda(\phi + 2\pi n)} e^{i\alpha\phi} I_{|\lambda|}\left(\frac{Mr r'}{2\hbar i\tau}\right), \end{aligned} \quad (20)$$

where $I_\nu(Z)$ is the modified Bessel function. In our case, we have $|\mathbf{r}'| = |P\mathbf{r}| = |\mathbf{r}|$, and $\phi = 0, \pi$:

$$K(\mathbf{r}, \mathbf{r}'; \tau) = \sum_{n=-\infty}^{\infty} \left(\frac{M}{4\pi\hbar i\tau} \right) \exp\left(\frac{-Mr^2}{2\hbar i\tau}\right) e^{in\phi} I_{|n-\alpha|}\left(\frac{Mr^2}{2\hbar i\tau}\right). \quad (21)$$

Therefore, we arrive at the result

$$\tilde{Z}_2 = \frac{1}{2} \sum_{\substack{n=-\infty \\ \text{even}}}^{\infty} \int_0^{\infty} dx e^{-x} I_{|n-\alpha|}(x). \quad (22)$$

This is, as expected, formally divergent. A convergence factor $e^{-\epsilon x}$ is inserted in the integrand, with $\epsilon \rightarrow 0$ at the end of the calculation. We use the result [15]

$$\begin{aligned} F_\nu(a) &\equiv \int_0^{\infty} dx e^{-ax} I_\nu(x) = \frac{1}{\sqrt{a^2 - 1}} (a + \sqrt{a^2 - 1})^{-\nu}, \\ F_\nu(1 + \epsilon) &\rightarrow \frac{1}{\sqrt{2\epsilon}} (1 + \sqrt{2\epsilon})^{-\nu}. \end{aligned} \quad (23)$$

As before, we appeal to eq. (8). Expanding about Fermi statistics, $\alpha = 2j + 1 + \delta$,

and $|n - \alpha| = 1 \pm \delta, 3 \pm \delta$, etc. Thus,

$$\begin{aligned} B(2j + 1 + \delta, T) &= \frac{1}{4}\lambda_T^2 + 2\lambda_T^2 \lim_{\epsilon \rightarrow 0} \\ &\left[\frac{1}{2} \frac{1}{\sqrt{2\epsilon}} \sum_{\substack{n=1 \\ \text{odd}}}^{\infty} (1 + \sqrt{2\epsilon})^{-n} (2 - (1 + \sqrt{2\epsilon})^\delta - (1 + \sqrt{2\epsilon})^{-\delta}) \right] \\ &= \frac{1}{4}\lambda_T^2 - \frac{1}{2}\delta^2\lambda_T^2. \end{aligned} \quad (24)$$

Expanding about Bose statistics introduces a term $|\delta|\lambda_T^2$ due to the $|n - \alpha| = |\delta|$ piece in the sum. Making the required periodic extension recovers the earlier result of eq. (12).

In some sense, the path integral result is more satisfying, because, although one still is presented with the delicacy involved with extracting the (finite) difference of two divergent expressions, there is no necessity to impose a finite volume constraint, which was originally effected in order to perform the mode counting. One might object to our original calculation on the grounds that the virial coefficient might possibly be sensitive to the manner in which we perform the mode counting, since the dominant terms in the sum of eq. (9) are those at the tail end. As we have seen, this fear is unfounded.

A striking result is the nonanalyticity of eq. (12). It would be interesting to know whether cusps also arise in higher-order virial coefficients.

Due to the proliferation of the number of relative angles, such higher-order virial coefficients are exceedingly difficult to evaluate. In the high density limit, one might consider averaging the statistical flux over the entire system, and then consider the effect of a net statistical uniform magnetic field of magnitude $B = n\alpha\phi_0$, where n is the particle density. It is possible to reproduce the correct form of the free energy to leading order in n in this manner, however, one loses periodicity in α , and only certain values of α actually yield the correct result.

The most significant feature of the statistical interaction is that it is long ranged, hence perturbation expansions in α yield divergences and resummation is necessary, a situation reminiscent of the electron gas. Nevertheless, this statistics transformation process does yield a viable method for interpolating quantum statistics. The representation of a Fermi gas in terms of a Bose gas may be useful in other contexts, such as lattice field theory. Unless the statistical interaction is treated nonperturbatively, however, divergences may be difficult to handle.

Finally, it is interesting to derive the lagrangian of eq. (14) for the solitons of the nonlinear σ -model. Let us briefly recall that the model in question has a unit-vector order parameter $n^a(x)$, $a = 1, 2, 3$ and a conserved topological current

$$J^\mu = \frac{1}{8\pi} \epsilon^{\mu\nu\lambda} \epsilon^{abc} n^a \partial_\nu n^b \partial_\lambda n^c. \quad (25)$$

The conservation of J^μ licenses us to manufacture a U(1) “gauge potential” by the curl equation

$$J^\mu = \epsilon^{\mu\nu\lambda} \partial_\nu A_\lambda. \quad (26)$$

The crucial point is that we could include a topological term

$$H = \frac{\Theta}{2\pi} \int d^3x A_\mu J^\mu \quad (27)$$

in the action, with Θ a real number ($\alpha \equiv \Theta/\pi$), which is analogous to the Θ -parameter in quantum chromodynamics. H is, in fact, the Hopf invariant describing maps of S^3 to S^2 . In a suitable gauge, such as $\partial A = 0$, we can solve for A_μ and so write H as a nonlocal interaction among the n^a fields. The solitons are bosons for $\Theta = 0$ ($\alpha = 0$) and fermions for $\Theta = \pi$ ($\alpha = 1$). In a more general context, any conserved current J_μ can be coupled to the vector field A_μ . If the only other appearance of A_μ in the lagrangian is the Chern-Simons term [16] $\epsilon_{\mu\nu\rho} A_\mu \partial_\nu A_\rho$, then A_μ represents a nondynamical field [17] which can be eliminated to give a nonlocal interaction, which will impart anomalous statistics to particles carrying charge associated with the current.

In ref. [3] the statistics of the solitons in this model were determined by invoking the linking number theorem. Here we will determine the statistics directly by interchanging two widely separated solitons, and in the process elucidate the linking number theorem.

For separations large compared to the sizes of the solitons we can approximate the solitons by point particles and the topological current by

$$J^\mu(x) = \sum_a \int d\tau \delta^{(3)}(x - q_a(\tau)) \frac{dq_a^\mu}{d\tau}, \quad (28)$$

with $a = 1, 2$ and $q_a(\tau)$ describing the trajectories of the two “point solitons.” We evaluate H by inserting eq. (28) into eqs. (26) and (27) and keeping only the cross-terms. The divergent self-interaction terms are evidently artifacts of the point approximation. To best understand the situation, we go to euclidean 3-space and think of eq. (26) as one of the time-independent Maxwell’s equations $\nabla \times \mathbf{B} = \mathbf{J}$ with the identification of A_μ as the magnetic field \mathbf{B} . Then H can clearly be interpreted as the work done on a magnetic monopole moving along the trajectory $q_1(\tau)$ by the magnetic field generated by an electric current flowing along the curve $q_2(\tau)$. With suitable normalization, this is just the number of times curves “1” and “2” wind around each other. We have thus made contact with the explicit form for the linking number between two curves given in mathematical texts [18]. This discussion also defines the linking number between two curves which are not closed.

To evaluate H explicitly, it is easiest to distort one of the curves, say “2,” to a straight line $q_2^\mu(\tau) = \tau \delta^{\mu 0}$, as we are allowed to do. We find by eq. (26) that

$A_i = \epsilon_{ij}x_j/r^2$, $A_0 = 0$, a pure (but topologically nontrivial) gauge. Once again, we could have appealed to $(2+1)$ -dimensional electrodynamics, this time interpreting J^0 as B . These remarks make clear that the effect here is essentially the Bohm-Aharonov phenomenon. It is sometimes convenient to transform to a singular gauge wherein $A = 0$ except along string singularities attached to each particle, across which A has a jump discontinuity of constant magnitude.

In summary, the action describing N of these point particles is just

$$S = \int dt L = \int dt \left[\frac{1}{2}m \sum_{a=1}^N \left(\frac{dx_a}{dt} \right)^2 + \frac{\Theta}{\pi} \sum_{a < b} \frac{d}{dt} \theta_{ab} \right]. \quad (29)$$

Here x_a is a two-dimensional vector locating particle a and θ_{ab} is the angle of particle b relative to particle a , measured from the x -axis, say. The preceding discussion has boiled ΘH down to the second term in this equation. As we have seen, although this term is a total time derivative and appears as an interaction, it determines the statistics of the particles.

In the original model, the solitons have topological charge $Q = \int d^3x J_0$ taking on all integer values. The $|Q| > 1$ solitons are unstable against breakup. In writing down eq. (29) we have included only $Q = +1$ particles. It is easy enough, however, to include $Q = -1$ particles as well by noting that the $(+-)$ “interaction” has opposite sign from the $(++)$ and $(--)$ “interactions.”

DPA would like to thank Stefan Theisen for making the work of Gerry and Singh known to us, and for many useful discussions. This work was supported in part by the National Science Foundation under grants DMR82-16285 and PHY77-27084, supplemented by funds from the National Aeronautics and Space Administration. One of us (DPA) is grateful for the support of an AT&T Bell Laboratories Scholarship.

Note added

This work supersedes the preprint “Interpolating quantum statistics”, NSF-ITP-84-25, by two of the authors (F.W. and A.Z.).

References

- [1] F. Wilczek, Phys. Rev. Lett. 49 (1982) 957
- [2] Y. Wu, US Dept. of Energy preprint 40048-09P4 (1984)
- [3] F. Wilczek and A. Zee, Phys. Rev. Lett. 51 (1983) 2250
- [4] R.B. Laughlin, Phys. Rev. Lett. 50 (1983) 1395
- [5] B.I. Halperin, Phys. Rev. Lett. 52 (1984) 1583
- [6] D. Arovas, R. Schrieffer and F. Wilczek, preprint NSF-ITP-84-66, submitted to Phys. Rev. Lett.
- [7] J.G. Dash, Films on solid surfaces (Academic Press, 1968)

- [8] M. Abramowitz and I. Stegun, *Handbook of mathematical functions* (Dover, 1972)
- [9] R.P. Feynman, *Statistical mechanics* (Benjamin, 1972)
- [10] Y. Aharonov and D. Bohm, *Phys. Rev.* 115 (1959) 485
- [11] C.C. Gerry and V.A. Singh, *Phys. Rev.* D20 (1979) 2550
- [12] A. Inomata and V.A. Singh, *J. Math. Phys.* 19 (1978) 2318
- [13] C.C. Gerry and V.A. Singh, *Nuovo Cim.* 73B (1983) 161
- [14] S.F. Edwards and Y.V. Gulyaev, *Proc. Roy. Soc. London* A279 (1964) 229
- [15] I.S. Gradshteyn and I.M. Ryzhik, *Table of integrals, series, and products* (Academic Pre
- [16] J. Schonfeld, *Nucl. Phys.* B185 (1981) 157;
S. Deser, R. Jackiw and S. Templeton, *Phys. Rev. Lett.* 48 (1982) 975;
Y. Wu, Washington preprint (1983)
- [17] C.R. Hagen, Rochester preprint (1983)
- [18] H. Flanders, *Differential forms* (Academic Press, 1963)

Chapter 6

THE QUANTIZED HALL EFFECT

- [6.1] D. Arovas, J. R. Schrieffer and F. Wilczek, “Fractional Statistics and the Quantum Hall Effect,” *Phys. Rev. Lett.* **53** (1984) 722–723 282
- [6.2] D. P. Arovas, “Topics in Fractional Statistics”* 284
- [6.3] S. M. Girvin and A. H. MacDonald, “Off-Diagonal Long-Range Order, Oblique Confinement, and the Fractional Quantum Hall Effect,” *Phys. Rev. Lett.* **58** (1987) 1252–1255 323
- [6.4] J. E. Avron, A. Raveh and B. Zur, “Quantum Conductance in Networks,” *Phys. Rev. Lett.* **58** (1987) 2110–2113 327
- [6.5] R. B. Laughlin, “Superconducting Ground State of Noninteracting Particles Obeying Fractional Statistics,” *Phys. Rev. Lett.* **60** (1988) 2677–2680 331
- [6.6] M. Wilkinson, “An Example of Phase Holonomy in WKB Theory,” *J. Phys.* **A17** (1984) 3459–3476 335

* Original Contribution.

6

The Quantized Hall Effect

The integrally quantized Hall effect¹ (IQHE) and the fractionally quantized Hall effect² (FQHE) were discovered in the very special, almost bizarre context of semiconductor heterostructures subjected to huge magnetic fields while held at millikelvin temperatures. However, the theories devised to describe the effects are quite novel and interesting, and it seems increasingly likely that some of the ideas that arise will find a much wider application.

Let us first discuss just what is observed. A layer of electrons may be trapped at the interface between two semiconductors, known as a heterojunction. The electrons in this layer can, to a good approximation, be idealized as a two-dimensional gas with Coulomb repulsion. The quantized Hall effect has to do with the behavior of a two-dimensional electron gas in a strong magnetic field, at low temperatures. It is found that under these conditions the Hall coefficient—that is, the ratio of current flow to transverse potential—behaves in a most peculiar way. To be specific, it is found that as the magnetic field is varied smoothly the Hall coefficient does not vary smoothly, but rather stays constant over finite intervals—“plateaus.” The plateaus are separated by intervals of more normal, continuous behavior.

Moreover, the value of the Hall coefficient on the plateaus is found to be

$$R = \frac{h}{\nu e^2} \quad (6.1)$$

where ν is an integer (for the IQHE) or a rational number with odd denominator (FQHE).* Of course, all integers are rational numbers with odd denominators, so the IQHE may be thought of as a special case of the FQHE. However, the most popular theoretical explanations of the two effects are quite different. Here we shall be most interested in the theory of the FQHE. (When we wish to refer to the two effects collectively, we shall speak of the QHE.)

* It is conventional and sometimes illuminating to display Planck's constant h in some of the equations. We shall however often lapse into using $\hbar = 1$ or $\hbar = 2\pi$ without warning.

Part of the interest of the quantized Hall effect is that the Hall coefficient is expressed in terms of fundamental physical quantities. This is quite remarkable, considering that the measurement is made directly on a macroscopic material, with all the complexity and “dirt” that implies. However, there is no question that the relation (6.1) is satisfied to a high degree of accuracy; for the integral Hall effect ($\nu = 1$) the equality has been established to better than a part in ten million.

To appreciate how unusual the QHE is, let us discuss it more quantitatively and contrast it with the more common behavior, which is essentially classical. The Lorentz force law reads

$$m \frac{d\mathbf{v}}{dt} = e \left(\mathbf{E} + \frac{\mathbf{v}}{c} \times \mathbf{B} \right), \quad (6.2)$$

and dissipative processes will lead to the Ohm's law behavior

$$\mathbf{j} = ne\mathbf{v} = \sigma_0 \left(\mathbf{E} + \frac{\mathbf{v}}{c} \times \mathbf{B} \right) = \sigma_0 \left(\mathbf{E} + \frac{1}{nec} \mathbf{j} \times \mathbf{B} \right) \quad (6.3)$$

Here, n is the electron density and σ_0 is the zero-field conductivity. We are concerned with a situation in which \mathbf{B} and \mathbf{E} are constant fields, with \mathbf{B} pointing perpendicular to the two-dimensional plane within which the electrons are confined while \mathbf{E} lies in that plane. We find then for the current the self-consistent equation

$$\begin{aligned} j_x &= \sigma_0 E_x + \frac{\sigma_0}{nec} j_y B \\ j_y &= \sigma_0 E_y - \frac{\sigma_0}{nec} j_x B \end{aligned} \quad (6.4)$$

or equivalently for the resistivity tensor ρ_{ij}

$$\begin{pmatrix} E_x \\ E_y \end{pmatrix} = \begin{pmatrix} \rho_{xx} & \rho_{xy} \\ \rho_{yx} & \rho_{yy} \end{pmatrix} \begin{pmatrix} j_x \\ j_y \end{pmatrix} = \begin{pmatrix} \rho_0 & B/nec \\ -B/nec & \rho_0 \end{pmatrix} \begin{pmatrix} j_x \\ j_y \end{pmatrix} \quad (6.5)$$

where $\rho_0 = \sigma_0^{-1}$. Notice especially the linear dependence of the Hall resistance ρ_{xy} on B . It is independent of σ_0 , reflecting its essentially non-dissipative character.

The basic observation of the quantized Hall effect is that as the density of electrons is varied at fixed magnetic field (or, more practically, if the magnetic field is varied at fixed electron density) the resistivity does not vary continuously, but rather is constant over finite intervals—the above-mentioned plateaus. Now comparing the FQHE and the classical values of the resistivity,

$$\frac{B}{nec} = \frac{h}{\nu e^2} \quad \text{or} \quad \frac{n}{B} = \frac{\nu e}{hc} \quad (6.6)$$

we see that the essence of the FQHE is that on the plateaus n/B is frozen at the rational value ν .

To interpret this further let us recall the basic quantum mechanics of charged particles in magnetic fields. As the textbooks teach us,³ the continuous spectrum of a free particle breaks up, in the presence of a magnetic field, into discretely spaced, highly degenerate levels known as Landau levels. The density of states in each Landau level is

$$n_0 = \frac{e}{hc} \quad (6.7)$$

per unit area (in units of the magnetic length squared), for each spin. Comparing to Eq.(6.6), we arrive at a simple interpretation of the fraction ν : it is the number of filled Landau levels. The FQHE reveals that certain filling fractions are especially stable, *i.e.*, energetically favorable; it is, in this sense, a commensurability effect.

(Many of the experiments correspond to densities such that only one Landau level is relevant, and furthermore such that the splitting between electron states having spin aligned or anti-aligned with the magnetic field is so great that only the former need be considered. For simplicity, we shall restrict ourselves to this case. Then the relevant electrons are all identical particles, and ν is less than one.)

With this introduction, we can intelligently begin to discuss the basic theory of the IQHE and FQHE. For those who wish to penetrate more deeply, we would recommend an extremely valuable, authoritative, and accessible account of the field as of 1986.⁴

The IQHE (that is, the case $\nu = 1$) is readily understood at the level of an independent particle model. In the presence of random impurities, the single-particle wave functions fall into two classes: a band of extended states and a set of spatially localized states. When the Fermi energy is in the region of localized states, varying the number of electrons only adds or subtracts localized states, which carry no current. Thus the magnitude of the current is stuck; and furthermore, it is stuck at the full Landau level value, as we can appreciate by turning on the impurities gradually from zero. Indeed, as long as no extended states cross the Fermi energy as the impurity potentials are turned on, the current remains at its value in the absence of impurities, *i.e.*, the full Landau level value. So the IQHE is an inevitable consequence of independent particle behavior, given a modicum of localization theory.

By implication then the FQHE must be an essentially collective, many-body effect. Great insight into its nature followed upon an inspired proposal, by Laughlin, of a variational wave function to describe it.⁵ The electron states in the lowest Landau level, in the symmetric gauge, are of the form

$$\psi(z) = f(z) e^{-\frac{1}{4\ell^2}|z|^2} \quad (6.8)$$

where $z = x + iy$, f is an arbitrary analytic function, and $\ell = (eB/c)^{1/2}$ is the magnetic length. Now the complete degeneracy of states in the Landau level means that there is effectively no kinetic energy associated with motion in the band. There is no inertia, and the particles are effectively infinitely massive. At the same time, they are subject to mutual Coulomb repulsion. To minimize the potential energy, one would like to space the electrons equally, forming a “Wigner crystal”.⁶ It might seem at first glance that there is no obstruction to doing just this, since it costs nothing to localize the electrons, the usual penalty imposed by the uncertainty principle being inoperative for infinitely massive particles. However, the restriction to the lowest Landau level makes it impossible to localize the electrons strictly. Their position is actually uncertain to within a magnetic length ℓ . At densities high enough so that the electrons at different sites of the putative Wigner crystal have overlapping wave functions, the crystal has every opportunity to melt. Nevertheless, we would expect that the most favorable many-body wave function is one that prevents the electrons from ever getting too close together, and that treats all the electrons symmetrically. Laughlin⁷ proposed a specific wave function with these desirable properties, *viz.*

$$\Psi^m(z_1, \dots, z_N) = \prod_{i < j} (z_i - z_j)^m e^{-\frac{1}{4\ell^2} \sum |z_i|^2} \quad (6.9)$$

up to normalization. Here, m must be an odd number so that Fermi statistics is respected. The Laughlin wave function is analytically tractable, and a number of its properties have been reliably deduced. For large numbers N of electrons, it represents electrons filling a disk of radius $N/2\pi\ell^2$, with density corresponding to filling fraction $1/m$. The zeroes of the wave function as particle coordinates coalesce effectively keeps the particles separated, and cuts down on the unfavorable Coulomb energy. Indeed the Laughlin wave function is known to be an exact ground state for certain model, short-range potentials, and there is good numerical evidence that it is an excellent approximation to the ground state for the actual Coulomb potential.⁸

Qualitatively, the Laughlin wave function represents an incompressible quantum liquid. Incompressibility, it will be seen, is the essence of the fractionally quantized Hall effect. Indeed we have interpreted the remarkable plateaus revealed by experiment as meaning that as n or B is varied, the filling fraction stays fixed. To account for this, we must find that the effect of adding a particle is not to change the filling fraction by a small amount over a large volume, but rather to leave the filling fraction pinned at its favorable value macroscopically, with the deviation from this value carefully localized. Thus, a density perturbation must lead to the creation of localized quasi-particles, with a gap in the spectrum corresponding to the finite energy cost of a quasi-particle. There are strong numerical indications that the quasi-particle (more accurately, quasi-hole) excitations around the Laughlin

ground state are well described by the wave function generated by acting on the ground state as follows:

$$\Psi_{z_0}^m \equiv \prod_i (z_i - z_0) \Psi^m. \quad (6.10)$$

This creates a quasi-hole at z_0 . Effectively, each electron feels a centrifugal barrier at z_0 , and is pushed away. Note that the wave function is little disturbed far away from z_0 , as required. Note too a certain resemblance between the polynomial factors in equations (6.9) and (6.10); after looking at these, the fact that quasiparticles have the quantum numbers of “fractional electrons” should seem quite plausible.

The Laughlin wave functions describe states at filling fractions $\nu = 1/m$, and indeed these are among the most prominent filling fractions observed experimentally in the FQHE. As the density moves away from a favored value, more and more quasi-particles are created. It is plausible that, when enough of them have been created, these quasi-particles in turn organize themselves into an incompressible quantum liquid, and so on. Implementing this thought, several authors have presented a hierarchical construction of states corresponding to other filling fractions.⁹ These constructions, although they certainly contain a core of truth, do not have the compelling simplicity and uniqueness of the Laughlin $1/m$ states.

Several of the enclosed papers make use of the geometric phase to elucidate important aspects of the FQHE. In their brief paper [6.1] Arovas, Schrieffer, and Wilczek derive the fractional charge and fractional statistics of quasiparticles in the FQHE by applying geometric phase techniques to the Laughlin wavefunction. The subsequent contribution [6.2] by Arovas, adapted from his thesis and published here for the first time, gives a much more detailed account of this calculation and other aspects of the quantum mechanics of fractional statistics particles, emphasizing the introduction of a gauge field to capture their dynamics, as we discussed in the previous chapter.

The Laughlin wave function, although it apparently yields an accurate estimate of the ground state energy, and incorporates the correct qualitative physics of the incompressible quantum liquid and its commensurability criterion, is perhaps not the end of all desire. In particular, it does not give a crisp answer to the question “What, precisely, characterizes the quantized Hall state?” comparable to the characterization of magnetic or superconducting states in terms of their order parameters. Also, the wave function is very specially adapted to the problem of electrons in strong magnetic fields, and it is not easy to generalize it to other situations where similar physics may be involved (see below).

In order to remedy these and other limitations, several physicists have attempted to define some kind of order parameter and an effective Landau-Ginzburg type theory for the FQHE.

By far the most substantial and interesting attempts along these lines have been made by Girvin and MacDonald [6.3], and by Read.¹⁰ These authors make essential use of the statistical gauge field that we met in Chapter 5.

There are two issues to address. One is the construction of an effective Lagrangian, that serves to summarize the low-energy excitations and their interactions in a compact fashion. A plausible idea, which is essentially what these authors propose and to some extent derive, is that the effective theory simply consists of fractionally charged particles with fractional statistics. We have seen that these are important qualitative properties of the FQHE quasiparticles, and it is not absurd to suppose that they are the only important long-range properties the effective theory need represent. Simple Lagrangians realizing these ideas are easy to construct, using the techniques discussed in the previous chapter and by Arovas below.

The second issue is the microscopic underpinning of the effective Lagrangian. Girvin and MacDonald have made important progress on this issue. The underlying idea, as we understand it, is to make a singular gauge transformation of the Laughlin wave function—which we know represents an anyon condensate—that removes the statistical phase factors, thus allowing the underlying condensate to reveal itself. Girvin and MacDonald show by explicit calculation that the transformed density matrix has (algebraic) long-range order.

A notable defect of the existing “effective field theories” of the FQHE is their failure to incorporate or to illuminate the commensurability effects. In these theories, the statistical parameter is given from outside, and there is no real understanding of why particular rational values are selected, nor (as far as we know) any connection to the hierarchical constructions. We believe that this failure will only be remedied by a deeper study of topological invariants of the effective gauge fields, that represent the statistical properties of the full wavefunction over the multi-particle configuration space.

As we have mentioned several times above, there are several other contexts in which one expects that insight gained from the quantized Hall effect will lead to progress. A comparatively simple but quite entertaining and possibly important one is the behavior of one-dimensional electron gases in complex topologies—i.e., networks of wires—in strong magnetic fields, as discussed in the enclosed paper of Avron [6.4].

More speculative but very exciting is an idea repeatedly mentioned by Anderson,¹¹ and given one concrete form by Kalmeyer and Laughlin,¹² that there is a qualitative resemblance between the physics of the FQHE and the

physics of several other notoriously difficult problems in condensed matter physics. An outstanding example is the Mott insulator problem. It has been realized since the 1930s that several metallic oxides are insulators even though they contain an odd number of electrons per unit cell. This contradicts the basic principles of band theory, since—because of electron spin—an odd number of electrons per unit cell should lead to a half-filled band, and thus to a metal. If the materials were antiferromagnetically ordered there would be no such difficulty, basically because the size of the unit cell would be doubled (at least). However, although several of the metallic oxides do exhibit antiferromagnetic order, others do not—and in fact in the paradigmatic case of NiO there is a Néel transition from an antiferromagnetic to a disordered state, but the material remains an insulator even above the transition. An attractive hypothesis is that the antiferromagnetic “spin solid” melts to a quantum spin liquid, just as in the FQHE the Wigner crystal melts to an incompressible liquid. The idea of a spin liquid is made more compelling by a variety of observations indicating that localized spin moments do persist through the Néel transition—the spins do not disappear, and in fact their short-range correlations seem to vary little through the transition; it is only the long-range order that disappears.

The Mott insulator problem has long been a skeleton in the closet of condensed matter theory. Due to the lack of a convincing theoretical framework in which to address it, and to the technological unimportance of the materials concerned, the problem has been widely ignored. However the recent discovery of high temperature superconductivity, in which anomalously insulating CuO insulators play a crucial role, has forced this problem back into the light of day. Kalmeyer and Laughlin devised a variational wave function of the FQHE type for a relevant model. Laughlin [6.5] subsequently pointed out that the quasi-particles around the ground state, defined by this variational wave function, carry fractional statistics, as should by now seem quite plausible. In fact, these so-called spinons are half-fermions—the phase i accumulates as one winds around another. Furthermore, it is plausible that holes injected by doping around the half-filled band bind to these spinons, making charged, half-fermionic quasiparticles called holons. Now half-fermions can be thought of as fermions with an additional attractive gauge interaction, and so it is not implausible that they should condense into a BCS type superconducting pair state. In fact, a composite particle composed of two half-fermions is easily seen to possess Bose statistics.

There is a distinct possibility that many-body effects could also occur, and be qualitatively important in electronic systems, well outside the thermodynamic limit, for example giving rise to commensurability effects in macromolecules. This idea is made more plausible by the fact that numerical simulations of the FQHE indicate pronounced energy minima at non-trivial rational fillings even for quite small systems.

The appearance of gauge fields in the effective description of FQHE states suggests that the relatively well-understood physics of the FQHE, including the notion of an incompressible quantum liquid, commensurability, and the associated rich phase structure, may have much to teach us about the even less-understood physics of gauge theories as studied in elementary particle physics. Let us note in particular that the phase structure of lattice quantum electrodynamics with a theta term exhibits some uncanny resemblances to the hierarchical structure of different FQHE states.¹³.

The final entry of this chapter, by Wilkinson [6.6], is more closely related to theories of the IQHE. It studies a system known as Harper's equation, which has been used to model the behavior of Bloch electrons in a perpendicular magnetic field. (The precise logical connection to the IQHE is a long story; we refer the reader to the article by Thouless in Ref. 4 for details.) An important feature of Harper's equation is that when the number of magnetic flux quanta per unit cell is rational, an energy gap appears. Wilkinson uses a WKB method to calculate the Bohr-Sommerfeld energy levels. As explained by Kuratsuji and Iida [3.9] in a more general context, to obtain the correct WKB quantization condition it is necessary to include the effect of Berry's phase. This turns out to be a considerably simpler way to derive the beautiful fractal spectrum of Harper's equation than previous methods.

-
- [1] K. v. Klitzing, G. Dorda, and M. Pepper, *Phys. Rev. Lett.* **45** (1980) 494.
 - [2] D.C. Tsui, H.L. Stormer, and A.C. Gossard, *Phys. Rev. Lett.* **48** (1982) 1559.
 - [3] L.D. Landau and E.M. Lifshitz, *Quantum Mechanics*, Course of Theoretical Physics vol. 3 (Oxford: Pergamon Press, 1977).
 - [4] R. Prange and S. Girvin, *The Quantum Hall Effect* (Berlin: Springer-Verlag, 1987).
 - [5] R.B. Laughlin, *Phys. Rev. Lett.* **50** (1983) 1395.
 - [6] E. Wigner, *Trans. Faraday Soc.* **34** (1938) 678.
 - [7] R.B. Laughlin, *Phys. Rev. Lett.* **50** (1983) 1395.
Chapter 7 in Ref. 4.
 - [8] F.D.M. Haldane, Chapter 8 in Ref. 4.
 - [9] F.D.M. Haldane, *Phys. Rev. Lett.* **51**, (1983) 605;
R. B. Laughlin, *Surface Science* **141**, (1984) 11;
B. Halperin, *Phys. Rev. Lett.* **52**, (1984) 1583.
 - [10] N. Read, unpublished.

- [11] P. Anderson, unpublished.
- [12] V. Kalmeyer and R.B. Laughlin, *Phys. Rev. Lett.* **59** (1987) 2095.
- [13] J. Cardy *Nucl. Phys.* **B205** (1982) 17;
A. Shapere and F. Wilczek, "Self-dual models with theta terms," to appear in *Nuclear Physics B*.

Fractional Statistics and the Quantum Hall Effect

Daniel Arovas

Department of Physics, University of California, Santa Barbara, California 93106

and

J. R. Schrieffer and Frank Wilczek

Department of Physics and Institute for Theoretical Physics, University of California, Santa Barbara, California 93106

(Received 18 May 1984)

The statistics of quasiparticles entering the quantum Hall effect are deduced from the adiabatic theorem. These excitations are found to obey fractional statistics, a result closely related to their fractional charge.

PACS numbers 73.40.Lq, 05.30.-d, 72.20.My

Extensive experimental studies have been carried out¹ on semiconducting heterostructures in the quantum limit $\omega_0\tau > 1$, where $\omega_0 = eB_0/m$ is the cyclotron frequency and τ is the electronic scattering time. It is found that as the chemical potential μ is varied, the Hall conductance $\sigma_{xy} = I_x/E_y = ve^2/h$ shows plateaus at $v = n/m$, where n and m are integers with m being odd. The ground state and excitations of a two-dimensional electron gas in a strong magnetic field B_0 have been studied²⁻⁴ in relation to these experiments and it has been found that the free energy shows cusps at filling factors $v = n/m$ of the Landau levels. These cusps correspond to the existence of an "incompressible quantum fluid" for given n/m and an energy gap for adding quasiparticles which form an interpenetrating fluid. This quasiparticle fluid in turn condenses to make a new incompressible fluid at the next larger value of n/m , etc.

The charge of the quasiparticles was discussed by Laughlin² by using an argument analogous to that used in deducing the fractional charge of solitons in one-dimensional conductors.⁵ He concluded for $v = 1/m$ that quasiholes and quasiparticles have charges $\pm e^* = \pm e/m$. For example, a quasi-hole is formed in the incompressible fluid by a two-dimensional bubble of a size such that $1/m$ of an electron is removed. Less clear, however, is the statistics which the quasiparticles satisfy; Fermi, Bose, and fractional statistics having all been proposed. In this Letter, we give a direct method for determining the charge and statistics of the quasiparticles.

In the symmetric gauge $\vec{A}(\vec{r}) = \frac{1}{2}\vec{B}_0 \times \vec{r}$ we consider the Laughlin ground state with filling factor $v = 1/m$,

$$\psi_m = \prod_{j < k} (z_j - z_k)^m \exp\left(-\frac{1}{4} \sum_i |z_i|^2\right). \quad (1)$$

where $z_j = x_j + iy_j$. A state having a quasi-hole localized at z_0 is given by

$$\psi_m^{+z_0} = N_+ \prod_i (z_i - z_0) \psi_m, \quad (2)$$

while a quasiparticle at z_0 is described by

$$\psi_m^{-z_0} = N_- \prod_i (\partial/\partial z_i - z_0/a_0^2) \psi_m, \quad (3)$$

where $2\pi a_0^2 B_0 = \phi_0 = hc/e$ is the flux quantum and N_\pm are normalizing factors.

To determine the quasiparticle charge e^* , we calculate the change of phase γ of $\psi_m^{+z_0}$ as z_0 adiabatically moves around a circle of radius R enclosing flux ϕ . To determine e^* , γ is set equal to the change of phase,

$$(e^*/\hbar c) \oint \vec{A} \cdot d\vec{l} = 2\pi(e^*/e)\phi/\phi_0, \quad (4)$$

that a quasiparticle of charge e^* would gain in moving around this loop. As emphasized recently by Berry⁶ and by Simon⁷ (see also Wilczek and Zee⁸ and Schiff⁹), given a Hamiltonian $H(z_0)$ which depends on a parameter z_0 , if z_0 slowly transverses a loop, then in addition to the usual phase $\int E(t') dt'$, where $E(t')$ is the adiabatic energy, an extra phase γ occurs in $\psi(t)$ which is independent of how slowly the path is traversed. $\gamma(t)$ satisfies

$$d\gamma(t)/dt = i \langle \psi(t) | d\psi(t)/dt \rangle. \quad (5)$$

From Eq. (2),

$$\frac{d\psi_m^{+z_0}}{dt} = N_+ \sum_i \frac{d}{dt} \ln(z_i - z_0(t)) \psi_m^{+z_0}, \quad (6)$$

so that

$$\frac{d\gamma}{dt} = iN_+^2 \left\langle \psi_m^{+z_0} \left| \frac{d}{dt} \sum_i \ln(z_i - z_0) \right| \psi_m^{+z_0} \right\rangle. \quad (7)$$

Since the one-electron density in the presence of

the quasihole is given by

$$\rho^{+z_0}(z) = \langle \psi_m^{+z_0} | \sum_i \delta(z_i - z) | \psi_m^{+z_0} \rangle, \quad (8)$$

we have

$$\frac{d\gamma}{dt} = i \int dx dy \rho^{+z_0}(z) \frac{d}{dt} \ln[z - z_0(t)], \quad (9)$$

where $z = x + iy$. We write $\rho^{+z_0}(z) = \rho_0 + \delta\rho^{+z_0}(z)$, with $\rho_0 = \nu B/\phi_0$. Concerning the ρ_0 term, if z_0 is integrated in a clockwise sense around a circle of radius R , values of $|z| < R$ contribute $2\pi i$ to the integral while $|z| > R$ contributes zero. Therefore, this contribution to γ is given by

$$\begin{aligned} \gamma_0 &= i \int_{|z| < R} dx dy \rho_0 2\pi i \\ &= -2\pi \langle n \rangle_R = -2\pi\nu\phi/\phi_0, \end{aligned} \quad (10)$$

where $\langle n \rangle_R$ is the mean number of electrons in a circle of radius R . Corrections from $\delta\rho$ vanish as $(a_0/R)^2$, where $a_0 = (\hbar c/eB)^{1/2}$ is the magnetic length. This term corresponds to the finite size of the hole.

Comparing with Eq. (4), we find $e^* = \nu e$, in agreement with Laughlin's result. A similar analysis shows that the charge of the quasiparticle $\psi_m^{-z_0}$ is $-e^*$.

To determine the statistics of the quasiparticles, we consider the state with quasiholes at z_a and z_b ,

$$\psi_m^{z_a, z_b} = N_{ab} \prod_i (z_i - z_a)(z_i - z_b) \psi_m. \quad (11)$$

As above, we adiabatically carry z_a around a closed loop of radius R . If z_b is outside the circle $|z_b| = R$ by a distance $d \gg a_0$, the above analysis for γ is unchanged, i.e., $\gamma = -2\pi\nu\phi/\phi_0$. If z_b is inside the loop with $|z_b| - R \ll -a_0$, the change of $\langle n \rangle_R$ is $-\nu$ and one finds the extra phase $\Delta\gamma = 2\pi\nu$. Therefore, when a quasiparticle adiabatically encircles another quasiparticle an extra "statistical phase"

$$\Delta\gamma = 2\pi\nu \quad (12)$$

is accumulated.¹⁰ For the case $\nu = 1$, $\Delta\gamma = 2\pi$, and the phase for interchanging quasiparticles is $\Delta\gamma/2 = \pi$ corresponding to Fermi statistics. For ν noninteger, $\Delta\gamma$ corresponds to fractional statistics, in agreement with the conclusion of Halperin.¹¹ Clearly, when ν is noninteger the change of phase $\Delta\gamma$ when a third quasiparticle is in the vicinity will depend on the adiabatic path taken by the quasiparticles as they are interchanged and the pair permutation definition used for Fermi and Bose statistics no longer suffices.

A convenient method for including the statistical phase $\Delta\gamma$ is by adding to the actual vector potential \mathbf{A}_0 a "statistical" vector potential $\tilde{\mathbf{A}}_\phi$ which has no independent dynamics. $\tilde{\mathbf{A}}_\phi$ is chosen such that

$$(e^*/\hbar c) \oint \tilde{\mathbf{A}}_\phi \cdot d\vec{l} = \Delta\gamma = 2\pi\nu, \quad (13)$$

when z_a encircles z_b . One finds this fictitious $\tilde{\mathbf{A}}_\phi$ to be

$$\tilde{\mathbf{A}}_\phi(\vec{r} - \vec{r}_b) = \frac{\phi_0 \hat{z} \times (\vec{r} - \vec{r}_b)}{2\pi |\vec{r} - \vec{r}_b|^2} \quad (14)$$

if the quasiparticles are treated as bosons and $\phi_0 \rightarrow \phi_0(1 - 1/\nu)$ if they are treated as fermions. Thus, the peculiar statistics can be replaced by a more complicated effective Lagrangian describing particles with conventional statistics.¹²

Finally, we note that if one pierces the plane with a physical flux tube of magnitude ϕ , the above arguments suggest that a charge $\nu e\phi/\phi_0$ is accumulated around the tube, regardless of whether ϕ/ϕ_0 is equal to the ratio of integers.

This work was supported in part by the National Science Foundation through Grant No. DMR82-16285 and No. PHY77-27084, supplemented by funds from the National Aeronautics and Space Administration. One of us (D.A.) is grateful for the support of an AT&T Bell Laboratories Scholarship.

¹K. von Klitzing, G. Dorda, and M. Pepper, Phys. Rev. Lett. **45**, 494 (1980).

²R. B. Laughlin, Phys. Rev. Lett. **50**, 1395 (1983).

³F. D. M. Haldane, Phys. Rev. Lett. **51**, 605 (1983).

⁴B. I. Halperin, Institute of Theoretical Physics, University of California, Santa Barbara, Report No. NSF-ITP-83-34, 1983 (to be published).

⁵W. P. Su and J. R. Schrieffer, Phys. Rev. Lett. **46**, 738 (1981).

⁶M. V. Berry, Proc. Roy. Soc. London, Ser. A **392**, 45-57 (1984).

⁷B. Simon, Phys. Rev. Lett. **51**, 2167 (1983).

⁸F. Wilczek and A. Zee, Phys. Rev. Lett. **52**, 2111 (1984).

⁹L. Schiff, *Quantum Mechanics* (McGraw-Hill, New York, 1955), p. 290.

¹⁰Although ψ is a variational wave function, rather than the actual adiabatic wave function, the statistical properties of the quasiparticles are not expected to be sensitive to this inconsistency. We could regard ψ to be an exact excited-state wave function for a model Hamiltonian.

¹¹B. I. Halperin, Phys. Rev. Lett. **52**, 1583, 2390(E) (1984).

¹²F. Wilczek and A. Zee, Institute of Theoretical Physics, University of California, Santa Barbara, Report No. NSF-ITP-84-25, 1984 (to be published).

Topics in Fractional Statistics

Daniel P. Arovas

Department of Physics
University of California at San Diego
La Jolla, CA 92093

ABSTRACT

I review the general theory of quantum mechanical particles with fractional statistics in two space dimensions ('anyons'). The thermodynamics of a low-density gas of anyons is discussed, as well as is the relevance of fractional statistics to quasiparticle excitations in the fractional quantized Hall effect.

Introduction

The physical implications of quantum statistics are numerous and lead to profound consequences for the excitation spectra and thermodynamic properties of all systems. Fermions, which obey the Pauli exclusion principle, fill a Fermi sphere in the $T \rightarrow 0$ limit.* Low temperature thermodynamics are then dominated by particle-hole excitations across the Fermi surface, leading to a specific heat $C_F^{d=3} \sim \frac{1}{2}\pi^2 Nk_B T/T_F$, which in the case of metals is much smaller than the Maxwell-Boltzmann value of $\frac{3}{2}Nk_B$ even at room temperature. Bosons, which do not heed the Pauli restriction of at most one quantum per single particle level i , exhibit a ground state in which the lowest such single particle level ($i = 0$) is macroscopically occupied; thermodynamic properties are then related to fluctuations in the occupancy of low-lying excited states. In dimensions $d > 2$, these fluctuations are weak enough to preserve the macroscopic occupancy of the $i = 0$ state, a phenomenon known as Bose condensation.

Even in the dilute gas regime, where the mean interparticle spacing, $n^{-1/2}$, is much larger than the thermal wavelength, $\lambda_T = (2\pi\hbar^2/mk_B T)^{1/2}$, relics of the quantum limit can be identified by examining corrections to the ideal gas law

$$p = nk_B T(1 + B_2 n + B_3 n^2 + \dots). \quad (1)$$

The terms B_i are the *virial coefficients* and characterize deviations from ideal gas behavior. For a free gas in two dimensions, $B_2(T) = \mp \frac{1}{4}\lambda_T^2$, the plus sign applying to the Fermi case, as the Pauli principle effectively pushes fermions away from each other, thus increasing the pressure.

Such well established properties ultimately derive from the respective symmetry and antisymmetry of Bose and Fermi wave functions, and when confronted with the query, “Why is there this sharp Bose–Fermi dichotomy; can no other quantum statistics be formulated?” most learned professors respond that indistinguishability implies that any N -body Hamiltonian will commute with elements σ of the permutation group S_N , and quantum mechanics, being a unitary theory, obliges us to characterize physical states by a one-dimensional representation of S_N , of which

* For the moment, I shall consider noninteracting particles, though the thermodynamic behavior discussed is qualitatively correct for interacting systems within the context of a quasiparticle model.

there are only two : the symmetric, or Bose representation ($\chi_B(\sigma) = 1$), and the antisymmetric, or Fermi representation ($\chi_F(\sigma) = \text{sgn}(\sigma)$). On the other hand, it is often emphasized that wave functions themselves are not physical entities; physical information is conveyed by propagators and matrix elements. It is then natural to ask whether the restrictions on quantum statistics discussed above also apply if one adopts a Feynman path integral approach to quantum mechanics, in which propagators, rather than wave functions, play the fundamental role.

A careful consideration of this issue leads to the possibility of exotic (or ‘fractional’) statistics in systems of spatial dimension less than three. In this chapter, I will discuss the case $d = 2$, for which a continuous one-parameter family of quantum statistics may be formulated — $d = 2$ also holds the possibility for physical relevance in the case of the fractional quantized Hall effect (FQHE). I will concentrate on the physical consequences of fractional statistics in the dilute gas regime and demonstrate explicitly how certain thermodynamic quantities interpolate between Bose and Fermi behavior as functions of the statistics determining parameter.

Charged Particle–Flux Tube Composites

The essential difference between systems of indistinguishable particles in two and three dimensions is easy to comprehend. In three or more spatial dimensions, no winding can be ascribed to relative particle motion because any two paths between a chosen pair of points in configuration space may be deformed into one another. Only when one descends below $d = 3$ does this relative winding become a well defined concept. It was realized by Wilczek¹ that this result could be exploited by associating to each particle a ‘charge’ e and a flux tube of strength $\phi = \alpha hc/e$ (see Fig.[1]). The coupling between the charges and the vector potential arising due to the flux tubes will then keep track of the relative winding of particles with an Aharonov–Bohm phase of $e^{-ie\phi/\hbar c}$ per revolution $\Delta\theta_{ij} = 2\pi$. These composites will obey fractional statistics — a feature which prompted Wilczek to name them ‘anyons’. The electromagnetic field due to the charges and fluxes is *non-dynamical*, *i.e.* its evolution is completely determined by the motion of the particles.

Consider the quantum mechanics of an electron confined to the two-dimensional plane with a flux tube of strength $\phi = \alpha\phi_0$ piercing the origin ($\phi_0 = hc/e$ is the Dirac flux quantum). The vector potential associated with the flux tube may be

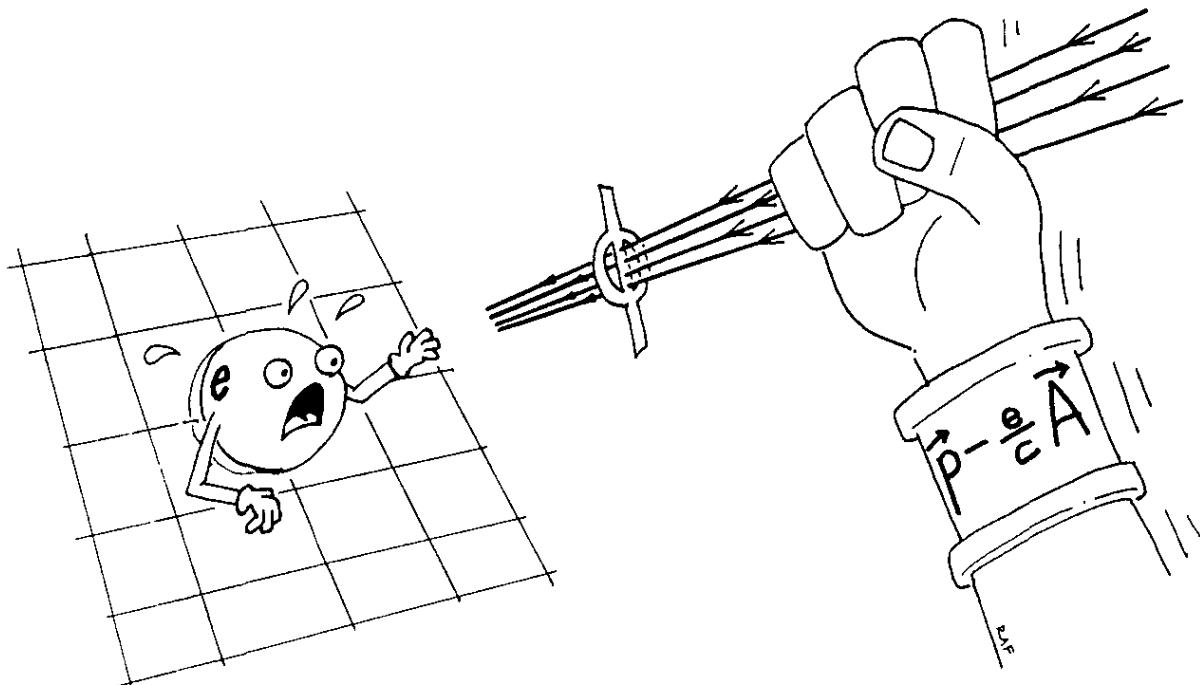


Figure 1. Artist's conception of a charged particle-flux tube composite.

taken as

$$\mathbf{A}(\mathbf{r}) = \frac{\alpha\phi_o}{2\pi} \frac{\hat{\mathbf{z}} \times \mathbf{r}}{r^2} = \frac{\alpha\phi_o}{2\pi} \nabla\theta. \quad (2)$$

The field strength, $\mathbf{B}(\mathbf{r}) = \nabla \times \mathbf{A} = \alpha\phi_o\delta(\mathbf{r})\hat{\mathbf{z}}$, is confined to the interior of the flux tube, which in this toy model is infinitesimally thin.

If the electron is otherwise free, the Hamiltonian is

$$H(\alpha) = \frac{1}{2m} \left(\mathbf{p} - \frac{e}{c} \mathbf{A} \right)^2 = -\frac{\hbar^2}{2m} \frac{1}{r} \frac{\partial}{\partial r} r \frac{\partial}{\partial r} + \frac{L_z^2(\alpha)}{2mr^2}, \quad (3)$$

where the operator $L(\alpha)$ is the familiar dynamical angular momentum,

$$L_z(\alpha) = e^{i\alpha\theta} \left(-i\hbar \frac{\partial}{\partial \theta} \right) e^{-i\alpha\theta} = \left(-i\hbar \frac{\partial}{\partial \theta} - \alpha\hbar \right). \quad (4)$$

If the system is placed in an eigenstate and the flux ϕ is varied adiabatically, the angular momentum (and hence the energy) will change with ϕ . Put simply, a variation in flux will lead to an azimuthal electric field $\mathbf{E}(\mathbf{r}) = -(\dot{\alpha}\hbar/e)r\hat{\theta}$ by Faraday's law.¹ The rate of change of angular momentum is then $\dot{L}_z = [\mathbf{r} \times (e\mathbf{E})]_z =$

$-\dot{\alpha}\hbar$, and therefore as the flux is cranked from $\phi = 0$ to $\phi = \alpha\phi_o$, the spectrum of allowed angular momenta changes from $\{L_z = m\}$ to $\{L_z = m - \alpha\}$.

An explicit solution to the eigenvalue equation $H(\alpha)\psi^\alpha = E(\alpha)\psi^\alpha$ is

$$\psi_k^\alpha(r, \theta) = \sum_{m=-\infty}^{\infty} A_m e^{im\theta} J_{|m-\alpha|}(kr) \quad (5)$$

where $J_\nu(x)$ is the Bessel function of the first kind of order ν . By imposing a hard wall constraint at $r = \Lambda$, $\psi_k^\alpha(\Lambda, \theta) = 0$, one can index the various modes by integers l and n :

$$\begin{aligned} \psi_{l,n}^\alpha(r, \theta) &= e^{in\theta} J_{|n-\alpha|}(k_{l,n}r) \\ k_{l,n} &= x_{|n-\alpha|, l}/\Lambda \\ E_{l,n} &= \hbar^2 k_{l,n}^2 / 2m \\ L_z |_{l,n} &= n - \alpha. \end{aligned} \quad (6)$$

(I have adopted the notation $x_{\nu,l}$ for the l^{th} node of J_ν .)

As expected, the energy spectrum is periodic in α , and an adiabatic increase of α by 1 results in the maps the spectrum back into itself by $\{l, n\} \rightarrow \{l - 1, n\}$. I would like to emphasize that Eq.(14) constitutes a gauge transformation only when α is an integer. Otherwise, the factors $e^{\pm ie\phi/\hbar c}$ are not single valued, and it is *incorrect*, however tempting, to write $\psi^\alpha = e^{i\alpha\theta}\psi^0$.

Two Anyons

I shall now discuss the quantum mechanics of two anyons. In the absence of an external potential, the Hamiltonian resembles $H = H_1 + H_2$, with

$$H_i = \frac{1}{2m} \left(\mathbf{p}_i - \frac{e}{c} \mathbf{A}_s(\mathbf{r}_i) \right)^2. \quad (7)$$

The ‘statistical vector potential’ $\mathbf{A}_s(\mathbf{r}_i)$ felt by particle 1 arises due to the presence of its companion and contains two identical contributions, one due to the interaction of the *charge* of 1 interacting with the *flux* of 2 and the other due to the interaction of the *flux* of 1 with the *charge* of 2. Thus,

$$\mathbf{A}_s(\mathbf{r}_1) = \frac{\alpha\phi_o}{2\pi} \frac{\hat{\mathbf{z}} \times (\mathbf{r}_1 - \mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|^2}, \quad (8)$$

where $\alpha = 2\phi/\phi_o$ is *twice* the number of Dirac flux quanta piercing any given particle; a corresponding expression applies for particle 2. By converting to relative

and center of mass coordinates (\mathbf{r}, \mathbf{R}) , one can easily solve this problem. In fact, the relative coordinate problem is just that of a particle and a flux tube, whose solution is displayed above. The two-body wave functions are then

$$\psi^\alpha(\mathbf{r}, \mathbf{R}) = e^{i\mathbf{K}\cdot\mathbf{R}} e^{in\theta} J_{|n-\alpha|}(kr). \quad (9)$$

The mode counting may be performed by introducing a potential $V(\mathbf{r})$ which forces the wave function to vanish when the interparticle separation exceeds some arbitrary value, Λ , which may later be taken to infinity. In this case, the allowed wavevectors are again quantized,

$$\begin{aligned} k_{l,n} &= x_{|n-\alpha|,l}/\Lambda \\ E_{\mathbf{K},l,n} &= \frac{\hbar^2 K^2}{4m} + \frac{\hbar^2 k_{l,n}^2}{m}, \end{aligned} \quad (10)$$

and all physical quantities are periodic in α .

If the particles themselves are taken to obey Bose statistics, the allowed values of n are restricted to the even integers. Similarly, Fermi statistics would require that n be odd. Therefore, according to Eq.(10) an adiabatic increase of α by 1 ($\phi \rightarrow \phi + \phi_0/2$) shifts the spectrum from Bose-like to Fermi-like and vice versa. More precisely, the eigenfunctions for a pair of fermions at $\alpha = 1$ are *unitarily equivalent* to those for a pair of bosons at $\alpha = 0$. Hence, no physical measurement could distinguish between the two cases and by the introduction of charge and flux, bosons can be magically transformed into fermions, a process which I shall refer to as ‘quantum alchemy’.

A General Recipe for Quantum Alchemy

Consider now an assembly of N particles interacting via arbitrary potentials. Perhaps there is also an external magnetic field $\mathbf{B}_{\text{ext}} = \nabla \times \mathbf{A}_{\text{ext}}$ present. To change the effective quantum statistics, one introduces the ‘statistical’ vector potential \mathbf{A}_s , arising from the flux tubes

$$\mathbf{A}_s(\mathbf{r}_i) = \frac{\alpha\phi_o}{2\pi} \sum_{j \neq i} \frac{\hat{\mathbf{z}} \times (\mathbf{r}_i - \mathbf{r}_j)}{|\mathbf{r}_i - \mathbf{r}_j|^2} \quad (11)$$

in which case the dynamical momenta become

$$\mathbf{p}_i^0 = \mathbf{p}_i - \frac{e}{c} \mathbf{A}_{\text{ext}}(\mathbf{r}_i) \rightarrow \mathbf{p}_i = \mathbf{p}_i - \frac{e}{c} \mathbf{A}_{\text{ext}}(\mathbf{r}_i) - \frac{e}{c} \mathbf{A}_s(\mathbf{r}_i). \quad (12)$$

The many-body Lagrangian one obtains is

$$L = \sum_i \left(\frac{1}{2} m \dot{\mathbf{r}}_i^2 - \frac{e}{c} \mathbf{A}_{\text{ext}}(\mathbf{r}_i) \cdot \dot{\mathbf{r}}_i \right) - V(\mathbf{r}_1, \dots, \mathbf{r}_N) - \alpha \hbar \frac{d}{dt} \sum_{i < j} \theta_{ij} \quad (13).$$

The Hamiltonian is then

$$H = \sum_{i=1} \frac{1}{2m} \left(\mathbf{p}_i - \frac{e}{c} \mathbf{A}_{\text{ext}}(\mathbf{r}_i) - \frac{e}{c} \mathbf{A}_s(\mathbf{r}_i) \right)^2 + V(\mathbf{r}_1, \dots, \mathbf{r}_N). \quad (14)$$

It is again worth emphasizing that although $H(\alpha) = e^{i\alpha\Theta} H(0) e^{-i\alpha\Theta}$, with $\Theta = \sum_{\text{pairs}} \theta_{ij}$, the function $\psi^\alpha = e^{i\alpha\Theta} \psi^0$ for general α is not properly single valued and hence does not constitute a solution to the Schrödinger equation with appropriate boundary conditions. The problem of obtaining single valued, N -body wave functions at arbitrary α from the $\alpha = 0$ solutions is in general an extremely difficult one. For $N = 2$, the problem is easily solved due to the separation of center of mass and relative coordinates. For $N = 3$, only a handful of solutions exist.² I know of no explicit wave functions describing fractional statistics for more than three particles. The situation gets worse with increasing particle number due to the fact that there are $\frac{1}{2}N(N-1)$ relative angles and only $(N-1)$ non-CM degrees of freedom. For $N = 2$, these numbers are the same.

General Theory

The laws for treating systems of indistinguishable particles by path integration were laid down by Laidlaw and DeWitt,³ who specialized to the case of $d = 3$. Following Wilczek's^{1,4} formulation of fractional statistics in terms of charged particle-flux tube composites, the general $d = 2$ theory was developed by Wu⁵ in 1984. Central to the formalism is the application of path integration in multiply connected spaces.⁶ The basic idea here is best described by Schulman's example of a free particle on a circle (S^1) parameterized by the angle ϕ . The propagator $K(\phi', t' | \phi, t)$ is written as a sum over paths

$$K(\phi', t' | \phi, t) = \sum_{n=-\infty}^{\infty} A_n \sum_{\phi(t) \in \mathcal{W}_n} e^{iS[\phi(t)]} \quad (15)$$

$$\mathcal{W}_n = \{\text{paths of winding number } n\}$$

where $S[\phi(t)]$ is the action corresponding to the path $\phi(t)$. The winding number of a path γ is given by the number of times γ passes some arbitrary point ϕ_0 ,

subtracting counterclockwise from clockwise passages. As Schulman points out, there is no *a priori* reason why the amplitudes A_n should be equal, as even with arbitrary A_n , the above expression correctly generates the Schrödinger equation when propagating a wave function over an infinitesimal time interval. Constraints on the A_n are dictated by requirements that the total probability $P = |K|^2$ be independent of ϕ_0 , and that K satisfy the composition rule

$$K(\phi'', t'' | \phi, t) = \int_0^{2\pi} d\phi' K(\phi'', t'' | \phi', t') K(\phi', t' | \phi, t). \quad (16)$$

This leads to the form $A_n = e^{i(\delta_0 + n\delta)}$.

The case of path integration on an arbitrary multiply connected manifold M follows by a straightforward generalization of the above discussion. The notion of winding number generalizes to that of homotopy, which is an equivalence relation used by topologists to classify paths[†] — two paths on M are homotopic if one can be continuously deformed into the other. The section below makes use of only the rudiments of homotopy theory, which are nicely summarized in Schulman's⁶ section 23.2, entitled "Rudiments of Homotopy Theory."

The propagator is written as a sum over homotopy classes

$$\begin{aligned} K(q', t' | q, t) &= \sum_{[\mu] \in \pi_1(M)} K_{[\mu]}(q', t' | q, t) \\ &= \sum_{[\mu] \in \pi_1(M)} \chi([\mu]) \sum_{q(t) \in [\mu]} e^{iS[q(t)]}, \end{aligned} \quad (17)$$

where $q \in M$ is a point, $[\mu] \in \pi_1(M)$ is a homotopy class, and $\pi_1(M)$ is the fundamental group of M . While homotopy among paths from q to q' is a perfectly well defined equivalence relation, it is convenient to label such paths by elements of the loop space of M . This is accomplished by defining a 'standard path mesh' consisting of a set of paths $C(q, q_0)$ from some arbitrary $q_0 \in M$ to every $q \in M$. This construction induces a mapping from $\pi_1(M, q, q')$ to $\pi_1(M, q_0)$:

$$\gamma: [0, 1] \mapsto M, \quad \gamma(0) = q, \quad \gamma(1) = q'$$

and

$$f_{qq'}: \pi_1(M, q, q') \mapsto \pi_1(M, q_0)$$

by

[†] More generally, a homotopy is an equivalence of maps between topological spaces.

$$f_{qq'}(\gamma) = C(q', q_0)^{-1} \circ \gamma \circ C(q, q_0). \quad (18)$$

Mesh invariance of the probability and the composition law require that χ be a unitary representation of $\pi_1(M)$. In the case of the circle,

$$\begin{aligned} M &= S^1, & \pi_1(S^1) &\cong \mathbf{Z} \\ \chi: \mathbf{Z} &\hookrightarrow U(1) & \text{by} & \quad \chi(n) = e^{i(\delta_0 + n\delta)}. \end{aligned} \quad (19)$$

The winding number on S^1 used above implicitly defined a path mesh consisting of paths which proceed directly from ϕ_0 to ϕ is a clockwise direction.

Because χ is a unitary representation, any nonabelian structure of $\pi_1(M)$ is left unused. Thus, any element of the commutator subgroup of π_1 will lie in the kernel of the homomorphism χ , and consequently it is only the *abelianized* π_1 (otherwise known as the first homology group H_1) which is needed. While the preceding statement sounds absolutely marvelous, its content is quite simple. Since $\chi([\mu])$ is a phase, both $|\nu|^{-1}[\mu]|\nu|$ and $[\mu]$ map to the same image under χ even though as elements of $\pi_1(M)$ they may be distinct. Roughly speaking, H_k measures the number of ' k -dimensional holes' in the manifold M , and physically, the association of a phase to each element of H_1 corresponds to threading each such hole with a magnetic flux ϕ .^f The phase accrued in winding about each hole is then $e^{-ie\phi/\hbar c}$, leading to a ϕ -dependent interference between paths of different winding number. Viewed as an element of S^1 , $e^{-ie\phi/\hbar c}$ distinguishes one among a continuum of quantum theories indexed by points on the circle.

Intuitively, one can think of this interference by taking a sort of cross product of M with its fundamental group. In this way, paths which start and end at the same point but exhibit different winding are not considered to be homotopic. This notion is but a barbarous interpretation of what algebraic topologists call a *universal covering space*. Every manifold M possesses a unique universal covering \tilde{M} and covering projection p with \tilde{M} simply connected and $p: \tilde{M} \hookrightarrow M$ a local homeomorphism. As discussed by Schulman, once one agrees on a specific choice of preimage p^{-1} for every point on M , any path on M may be lifted to \tilde{M} . Consider two paths in \tilde{M} each emanating from \tilde{x} and concluding at distinct points \tilde{y}_a and \tilde{y}_b , respectively. If $p(\tilde{x}) = x$ and $p(\tilde{y}_a) = p(\tilde{y}_b) = y$, these paths may be viewed as the lifts of

^f Technically, one may associate a different flux ϕ_a with each generator a of H_1 . Think about puncturing the plane at n distinct points.

two nonhomotopic paths in M . In this way, the restricted propagator $K_{[\mu]}(q', t' | q, t)$ may be calculated by summing over *all* paths on the universal covering space from \tilde{q} to a particular $\tilde{q}_{[\mu]}$.

These concepts are nicely illustrated by the simple example of the 2-torus, $\mathbf{T}^2 = \mathbf{S}^1 \times \mathbf{S}^1$, shown in Fig.[2]. As one would expect, $\pi_1(\mathbf{T}^2) \cong H_1(\mathbf{T}^2) \cong \mathbf{Z} \times \mathbf{Z}$, and the universal covering space is the plane, \mathbf{R}^2 . The canonical projection p is defined as

$$p: \mathbf{R}^2 \rightarrow \mathbf{T}^2 \quad \text{by} \quad p(x, y) = (e^{ix}, e^{iy}). \quad (20)$$

The path $\{\tilde{\gamma}: [0, 1] \rightarrow \mathbf{R}^2 \mid \tilde{\gamma}(t) = (8\pi t, -2\pi t)\}$ is the lift of a path $\gamma: [0, 1] \rightarrow \mathbf{T}^2$ with winding numbers $n_1 = 4$ and $n_2 = -1$. In general, the phase $\chi(\gamma)$ may be defined in terms of the lifted path $\tilde{\gamma} = p^{-1}(\gamma)$:

$$\begin{aligned} \chi(\gamma) &= \exp \left(-i \frac{e}{\hbar c} \int_{\tilde{\gamma}} \mathbf{A} \cdot d\mathbf{l} \right) \\ \mathbf{A} &= (\phi_1, \phi_2). \end{aligned} \quad (21)$$

Notice that in the above equation, $\chi(\gamma)$ depends continuously on the endpoints, which seems to violate the requirement that χ be a homomorphism of the *discrete* group $\pi_1(M)$ onto $U(1)$. Indeed, by judicious use of the path mesh, this apparent inconsistency could be eliminated. But there really is no problem at all, for the respective propagators $K(q', t' | q, t)$ will differ only by a $U(1)$ -valued function of q_0 (the arbitrary base point), q , and q' . Since this will not affect the probability $P = |K|^2$, it is permissible to redefine the propagator with this phase included.

Configuration Space

Consider now a set of N particles which are confined to a manifold M . M may or may not be simply connected, and $d = \dim(M)$ is as yet unspecified, though the easiest case to imagine is of course $M = \mathbf{R}^d$. For distinguishable particles, the N -particle configuration space is simply the N -fold product $\mathcal{D}_N(M) \equiv M \times M \times \dots \times M$. For indistinguishable particles, one might expect the appropriate space to be $\mathcal{D}_N(M)/S_N$. Unfortunately, such a space is not a manifold. The problem lies in the coincidence points, *i.e.* the set $\mathcal{C}_N(M) \equiv \{(p_1, \dots, p_N) \in \mathcal{D}_N(M) \mid p_j = p_k \text{ for some } j, k\}.$ * To see this, consider the simpler case of two indistinguishable

* In general, when one takes the quotient of any manifold M with a discrete group G , the resulting space will not be a manifold if M contains points which are fixed under the action of elements of G (aside from the identity).

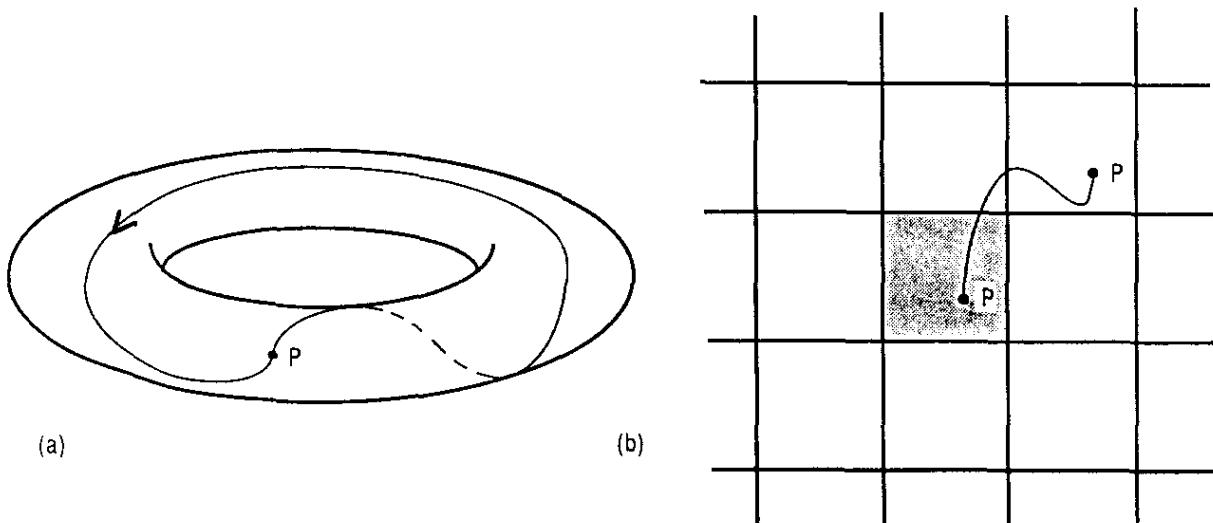


Figure 2. (a) The path shown on the torus is noncontractible, i.e. it cannot be continuously deformed to a single point. (b) When the curve is lifted to the universal covering space of the torus (the plane), the end points lie in different principal regions. Any such curve may be characterized by an element (n_1, n_2) of the fundamental group of the torus, $\pi_1(T^2) \cong \mathbf{Z}^2$, where n_1 and n_2 measure the total number of boundary crossings between squares in the horizontal and vertical directions, respectively.

particles on a line. The candidate configuration space is then $\mathbf{R} \times \mathbf{R}/S_2$, which is simply the set of points in \mathbf{R}^2 on or below the diagonal. The boundary imposed by the diagonal ruins the manifold structure. To cope with this mathematical nuisance, Laidlaw and DeWitt proposed that all coincidence points be excluded from $\mathcal{D}_N(M)$ before ‘modding out’ by S_N . The resulting configuration space is then $I_N(M) \equiv (\mathcal{D}_N(M) - \mathcal{C}_N(M)) / S_N$. Physically, this situation would be realized only if the particles had infinitely hard cores, and although this doesn’t seem to affect the generality of the discussion in any essential way, I know of no proof of this.

It is the space $I_N(M)$ which enters into the many particle path integral. Paths in this space may be classified by elements of the loop space $\mathcal{B}_N(M) \equiv \pi_1(I_N(M))$, as previously discussed. The structure of the loop space will determine what are the allowed phase factors $\chi([\mu])$ which multiply the restricted propagators $K_{[\mu]}$ in the sum of Eq.(17).

For $d > 2$, $\mathcal{D}_N(M)$ is simply connected, and $\mathcal{B}_N(M)$ is isomorphic to S_N , the only unitary representations of which are χ_B and χ_F . Thus, only Bose and Fermi statistics are allowed for spatial dimension three or higher. This was the main

conclusion of Laidlaw and DeWitt. In two dimensions, however, $\mathcal{D}_N(M)$ is multiply connected, and $\mathcal{B}_N(M)$ has the structure of an infinite nonabelian group, known in mathematical parlance as the full N -string braid group on M .

The structure of the braid groups for several two-dimensional manifolds has been investigated by Fadell and Van Buskirk,⁷ Birman,⁸ and others.⁹ For the relevant case of $M = \mathbf{R}^2$, there are $N - 1$ generators σ_i of $\mathcal{B}_N(\mathbf{R}^2)$ which obey the relations

$$\begin{aligned}\sigma_j \sigma_k &= \sigma_k \sigma_j, & |k - j| &\geq 2 \\ \sigma_i \sigma_{i+1} \sigma_i &= \sigma_{i+1} \sigma_i \sigma_{i+1}, & i &= 1, \dots, N - 2\end{aligned}\tag{22}$$

from which it is clear that any representation χ must satisfy $\chi(\sigma_i) = \chi(\sigma_{i+1})$ for all i , i.e. all generators σ_i map to *the same* unitary phase $e^{-i\pi\alpha}$ under the homomorphism χ . Therefore, there exists a continuous one-parameter family of unitary representations of $\mathcal{B}_N(\mathbf{R}^2)$ indexed by α .

There is a simple physical interpretation, due to Wu, of the algebra of Eq.(22). Consider a path $\gamma: [0, 1] \mapsto \mathcal{I}_N(\mathbf{R}^2)$. This path corresponds to a set of N ‘world-lines’ in the three-dimensional slab $[0, 1] \times \mathbf{R}^2$. These world-lines will in general interweave and be tangled like strings; elements of the braid group will describe the nature of this web. The association of an element of $\mathcal{B}_N(\mathbf{R}^2)$ with the path γ proceeds in several steps. First, obtain a labeling $1, \dots, N$ for the particles by projecting them onto the x -axis. Thus, the label 1 applies to that particle whose abscissa is the smallest. This labeling scheme is possible at all times t during the trajectory, except for those moments when two particles have abscissae which coincide. Such events are termed *crossings*; it is obvious that only particles with consecutive labels can cross. To each crossing one associates a factor $\sigma_n^{\pm 1}(t_x)$, where n and $n + 1$ participate in the crossing at time t_x , and the plus sign is used if the ordinate for n is greater than the ordinate for $n + 1$ at the crossing. The resulting product is then time ordered, with earlier time arguments appearing to the right. Finally, the time labels are discarded, and one is left with a product of σ_i ’s and σ_j^{-1} ’s, that is to say, an element of $\mathcal{B}_N(\mathbf{R}^2)$.

As discussed previously, not all of the detailed structure of $\pi_1(\mathcal{B}_N(\mathbf{R}^2))$ is needed in defining the propagator. The phase $\chi(\gamma)$ will register a factor $e^{-i\pi\alpha}$ for each crossing, and hence it is only the ‘net winding’ of all the particles, measured by the sum of the angular differences $\sum_{\text{pairs}} \Delta\theta_{ij} \equiv \sum_{\text{pairs}} \int dt \dot{\theta}_{ij}$, to which $\chi(\gamma)$

is sensitive.[†]

I will now describe the general form for the propagator on $\mathcal{I}_N(\mathbf{R}^2)$. Appealing to the earlier example of the torus, the propagator $K(q', t' | q, t)$ between $q = \{\mathbf{r}_1, \dots, \mathbf{r}_N\}$ and $q' = \{\mathbf{r}'_1, \dots, \mathbf{r}'_{N'}\}$ is expressed as a sum over all paths $\tilde{q}(t)$ in the universal covering space $\tilde{\mathcal{I}}_N(\mathbf{R}^2)$ subject to $p(\tilde{q}(t_k)) = q_k$ for $k = 1, 2$. The space $\tilde{\mathcal{I}}_N(\mathbf{R}^2)$ is simply $\mathcal{D}_N(\mathbf{R}^2)$ with the domains of each of the $N(N-1)/2$ relative angles θ_{ij} extended from the interval $[0, 2\pi]$ to the entire real line. The canonical projection p then maps each of the θ_{ij} back onto $[0, 2\pi]$ by taking the remainder upon dividing by 2π ; it also respects the permutation symmetry. Therefore,

$$K(q', t' | q, t) = \sum_{\sigma \in S_N} \sum_{|\mu|} \chi(|\mu|) K_{|\mu|}(\sigma q', t' | q, t) \quad (23)$$

$$\sigma q' \equiv \{\mathbf{r}'_{\sigma(1)}, \dots, \mathbf{r}'_{\sigma(N)}\}.$$

If the dynamics are determined by a Lagrangian $L(q, \dot{q})$, the resulting path integral is given by

$$K^\alpha(q', t' | q, t) = \sum_{\sigma \in S_N} \int \mathcal{D}q(\tau) \exp \left[i \int_q^{\sigma q'} d\tau \left(L - \alpha \frac{d}{d\tau} \sum_{i < j} \theta_{ij}(\tau) \right) \right]. \quad (24)$$

It is clear that the parameter α determines the statistics of the particles. Consider the diagonal elements of K^α for which $q = q'$. With $\alpha = 0$, all permutations carry the same signature and the Bose propagator is recovered. When $\alpha = 1$, the phase χ is given by the Fermi factor

$$(-1)^{\{\# \text{ of pairwise exchanges}\}}.$$

(Note also that the probability $P = |K|^2$ is periodic under $\alpha \rightarrow \alpha + 2$.) For noninteger α , K^α interpolates between these two familiar cases, and the statistics are said to be *fractional*. Note also that the Lagrangian defined by Eq.(24) is identical to that derived earlier in Eq.(13).

It should be stressed that α couples to a boundary term, namely the net winding number, and consequently it does not appear in the equations of motion.

[†] Since θ_{ij} is not single valued, $d\theta_{ij}/dt$ cannot be regarded as an exact differential. The idea here is that there is a distinction between $\Delta\theta_{ij} = 0$ and $\Delta\theta_{ij} = 6\pi$, for example.

Anyon Thermodynamics

All thermodynamic properties I will discuss derive from the behavior of the grand partition function,

$$\begin{aligned}\Xi(T, A, \mu) &= e^{-\beta\Omega} = \text{Tr } e^{-\beta(H - \mu N)} \\ &= \sum_{N=0}^{\infty} z^N Z_N,\end{aligned}\tag{25}$$

where μ is the chemical potential, $\beta = 1/k_B T$ the inverse temperature, and A the area of the system. Adhering to convention, I shall abbreviate $z = e^{\beta\mu}$ as the fugacity and Z_N as the N -particle (ordinary) partition function. Partial derivatives of the grand potential Ω with respect to its arguments yield entropy, pressure, and particle number:

$$\begin{aligned}d\Omega &= -SdT - pdA - Nd\mu \\ S &= -\frac{\partial\Omega}{\partial T}\Big|_{A, \mu} \quad p = -\frac{\partial\Omega}{\partial A}\Big|_{\mu, T} \quad N = -\frac{\partial\Omega}{\partial\mu}\Big|_{T, A}.\end{aligned}\tag{26}$$

In order to represent the pressure as a function $p(T, A, N)$, as in Eq.(1), one must turn a thermodynamic crank which inverts the defining relations of Eq.(26):

$$\begin{aligned}p/k_B T &= -\Omega/Ak_B T = \sum_{l=1}^{\infty} b_l(T, A) z^l \\ n &= N/A = \sum_{l=1}^{\infty} l b_l(T, A) z^l \\ \implies p/k_B T &= \sum_{l=1}^{\infty} B_l(T, A) n^l\end{aligned}\tag{27}$$

The virial coefficients $B_l(T, A)$ are so obtained. For classical interacting systems, the coefficients $b_l(T, A)$ are the well known Mayer cluster integrals. In general, the b_l may be expressed in terms of the first l many-body partition functions. For example, the second virial coefficient, which determines the lowest order corrections to ideal gas behavior, is given by

$$B_2(T, A) = A \left(\frac{1}{2} - Z_2/Z_1^2 \right).\tag{28}$$

The thermodynamics of the ideal quantum gases is discussed in many textbooks. It is essentially an academic exercise to work out the low density quantum corrections to Maxwell-Boltzmann statistics — an exercise rendered more academic by

the fact that no such systems have ever been observed in nature. Perhaps the closest that workers have come to isolating such quantum corrections in two-dimensional systems is in studies of ^4He and ^3He films adsorbed on Grafoil.¹⁰ Even in these systems, however, it is the interparticle potential *in conjunction* with the statistics that determines the thermodynamic behavior, rather than pure quantum corrections alone. For example, the measured heat capacity exhibits small deviations from the Dulong-Petit law $C = Nk_B$. A truncated virial expansion would read

$$C/Nk_B = 1 - n\beta^2 \frac{d^2 B_2}{d\beta^2} \quad (29)$$

and, since the ideal quantum gas would give $B_2 = \mp \frac{1}{4}\lambda_T^2 \propto \beta$, there would be *no* correction to C at this order. The helium atoms are, of course, interacting, and experience a hard core repulsion together with a weakly attractive van der Waals tail. As discussed by Dash,¹¹ a ^4He pair, whose ground state relative coordinate wave function is an *s*-wave, is more able to sample the repulsive core than a ^3He pair, which obeys the exclusion principle, leading to qualitatively different physics. With the interparticle potential properly treated, the second virial coefficient will indeed give a contribution to Eq.(29).

I shall be concerned with *purely* quantum effects in the thermodynamics of particles obeying fractional statistics. Since the statistical interaction is a many-body one (see Eq.(24) above), the problem is highly nontrivial. In fact, I shall only discuss the second virial coefficient, as higher order terms in the virial expansion require detailed knowledge of three-anyon dynamics, which is not yet available.

The thermodynamic functions for two-dimensional ideal quantum gases may be calculated from the following formulae:

$$\begin{aligned} \Omega(T, A, z) &= \mp Ak_B T \lambda_T^{-2} \zeta_2(\pm z) \\ n(T, z) &= \mp \lambda_T^{-2} \ln(1 \mp z) \\ F(n, T) &= \Omega + \mu N = \mp Ak_B T \lambda_T^{-2} \zeta_2(1 - e^{\mp n \lambda_T^2}) + n Ak_B T \ln \left(\mp(e^{\mp n \lambda_T^2} - 1) \right) \end{aligned} \quad (30)$$

where I write the generalized Riemann zeta function

$$\zeta_p(z) \equiv \sum_{n=1}^{\infty} \frac{z^n}{n^p}. \quad (31)$$

(In all formulae such as Eq.(30), the top sign shall refer to bosons, and the bottom sign to fermions.) The dimensionless parameter which characterizes these functions

is the scaled density $\varrho \equiv n\lambda_T^2$, which is the mean number of particles inside a square thermal wavelength. In the quantum limit $\varrho \rightarrow \infty$, F tends to the ground state energy, which is finite for fermions and zero for bosons. The Maxwell-Boltzmann limit ($\varrho \rightarrow 0$) is one of low density and/or high temperature, and is thus entropy dominated, the free energy tending to negative infinity as

$$F/Nk_B T = \ln(\varrho) - 1 + \sum_{k=1}^{\infty} \frac{1}{k} (B_{k+1} \lambda_T^{-2k}) \varrho^k. \quad (32)$$

The virial expansion converges well in the Maxwell Boltzmann limit. [†]

It would be nice to have explicit formulae for general α that could interpolate between the Fermi and Bose behavior described by Eq.(30). One would then witness a spectacular change in low temperature properties as the ground state changes from a filled Fermi disk to a Bose condensate. Even with no interparticle potential, I unfortunately see little prospect for such a complete solution, due to the extreme nastiness of the general many-anyon problem. Here I will describe a calculation of $B_2(T)$ which itself evidences some peculiar nonanalytic structure.*

Anyons in the Low Density Limit

My aim here is to calculate $B_2(\alpha, T)$ for the anyon gas. According to Eq.(28), this requires an evaluation of the two-particle partition function, Z_2 — this is possible due to the separability of H_2 into H_{rel} and H_{CM} . For the purpose of increased generality, I shall solve this problem in the presence of a constant magnetic field of arbitrary strength. The two-body Hamiltonian is written as a sum

$$\begin{aligned} H_2 &= H_{\text{CM}} + H_{\text{rel}} \\ H_{\text{CM}} &= \frac{1}{2M} \left(\mathbf{P} + \frac{1}{2} M \omega_c \hat{\mathbf{z}} \times \mathbf{R} \right)^2 \\ H_{\text{rel}} &= \frac{1}{2\mu} \left(\mathbf{p} + \frac{1}{2} \mu \omega_c \hat{\mathbf{z}} \times \mathbf{r} + \alpha \hbar \frac{\hat{\mathbf{z}} \times \mathbf{r}}{r^2} \right)^2 \end{aligned} \quad (33)$$

[†] It is a straightforward task to derive $\tilde{B}_3 = -1/36$, $\tilde{B}_4 = 0$, $\tilde{B}_5 = -1/3600$, $\tilde{B}_6 = 0$, $\tilde{B}_7 = 1/211680$, $\tilde{B}_8 = 0$, $\tilde{B}_9 = 1/10886400$, etc, for $\tilde{B}_n \equiv B_n \lambda_T^{2-2n}$. These results apply to both the Fermi and Bose cases. I thank MACSYMA for these numbers.

* In the thermodynamic limit, one defines $B_2(T) = \lim_{A \rightarrow \infty} B_2(T, A)$.

where $M = 2m$ is the total mass, $\mu = \frac{1}{2}m$ is the reduced mass (not to be confused with the chemical potential), and $\omega_c = eB/mc$ is the cyclotron energy. The general expression for the propagator, Eq.(24), is now continued to imaginary time:

$$Z_2 = \text{Tr } e^{-\beta H_2} = \frac{1}{2} \int d^2R d^2r \left[\langle \mathbf{R}, \mathbf{r} | e^{-\beta H_2} | \mathbf{R}, \mathbf{r} \rangle + \langle \mathbf{R}, -\mathbf{r} | e^{-\beta H_2} | \mathbf{R}, \mathbf{r} \rangle \right] \quad (34)$$

The center of mass contribution thereby factors out, yielding

$$Z_{\text{CM}} = 2Z_1 = 2A\lambda_T^{-2} \left(\frac{1}{2}\beta\hbar\omega_c \right) \text{csch} \left(\frac{1}{2}\beta\hbar\omega_c \right). \quad (35)$$

To obtain Z_{rel} , one needs the propagator $K^\alpha(\mathbf{r}', \mathbf{r}''; \tau)$ for a single particle in the presence of a flux tube. All the formalism necessary to tackle this problem was derived by Edwards and Gulyaev,¹² Peak and Inomata,¹³ and others¹⁴ in addressing the issue of path integration in polar coordinates, a somewhat subtle affair. A brief, but complete review is given in appendix A.

According to Eqs.(A.21-A.24), the imaginary time relative propagator is

$$\begin{aligned} K(\mathbf{r}', \mathbf{r}''; -i\beta\hbar) &= \frac{\mu\omega_c}{4\pi\hbar} \text{csch} \left(\frac{1}{2}\beta\hbar\omega_c \right) \\ &\times \exp \left[-\frac{\mu\omega_c}{4\hbar} \text{ctnh} \left(\frac{1}{2}\beta\hbar\omega_c \right) (r'^2 + r''^2) \right] e^{-\frac{1}{2}\beta\hbar\omega_c\alpha} \\ &\times \sum_{n=-\infty}^{\infty} e^{-\frac{1}{2}\beta\hbar\omega_c n} e^{in(\theta'' - \theta')} I_{|n+\alpha|} \left[\frac{\mu\omega_c}{2\hbar} r' r'' \text{csch} \left(\frac{1}{2}\beta\hbar\omega_c \right) \right], \end{aligned} \quad (36)$$

where $I_\nu(x)$ is the modified Bessel function of the first kind of order ν . I now define

$$\begin{aligned} F_+(-i\beta\hbar) &\equiv \int d^2r K(\mathbf{r}, \pm\mathbf{r}; -i\beta\hbar) \\ &= \frac{1}{2} e^{-\alpha\Delta} \sum_{n=-\infty}^{\infty} (\pm 1)^n e^{-n\Delta} \int_0^\infty dx e^{-x \cosh \Delta} I_{|n+\alpha|}(x) \end{aligned} \quad (37)$$

and

$$\Delta \equiv \frac{1}{2}\beta\hbar\omega_c. \quad (38)$$

Notice that F_+ is periodic in α with period 2. The partition function Z_{rel} may be written

$$Z_{\text{rel}} = \frac{1}{2} [F_+(-i\beta\hbar) \pm F_-(-i\beta\hbar)] \quad (39)$$

with the top sign applying when the fiducial ($\alpha = 0$) statistics refer to bosons. Integrals such as those arising in Eq.(37) may be evaluated exactly:¹⁵

$$\int_0^\infty dx e^{-ax} I_\nu(x) = \frac{1}{\sqrt{a^2 - 1}} \left(a + \sqrt{a^2 - 1} \right)^{-\nu}. \quad (40)$$

Since the partition function scales with the size of the system, one expects the sum in Eq.(37) to diverge. Indeed this is so. Consequently, in evaluating the second virial coefficient via Eq.(38), one is presented with the delicate task of extracting the finite difference of two divergent expressions (namely $\frac{1}{2}A$ and AZ_2/Z_1^2). In order to accomplish this, I shall employ a simple regularization procedure which renders finite all integrals such as those encountered in Eq.(37). This amounts to introducing a factor

$$\exp\left(-\epsilon \frac{\mu\omega_c r^2}{2\hbar} \operatorname{csch}\Delta\right) \quad (41)$$

under every integral $\int d^2r$. The subtraction shall then be performed on the regularized quantities, with $\epsilon \rightarrow 0$ at the end of the calculation.

To demonstrate how the regularization procedure works, I shall first evaluate $B_2(\alpha = 0, T)$ within this scheme. In this case, the identity

$$\sum_{n=-\infty}^{\infty} (\pm 1)^n e^{-n\Delta} I_{|n|}(x) = e^{\pm x \cosh \Delta} \quad (42)$$

gives rise to the simple expressions

$$\begin{aligned} F_+(-i\beta\hbar) &= \frac{1}{2} \int_0^\infty dx \\ F_-(-i\beta\hbar) &= \frac{1}{2} \int_0^\infty dx e^{-2x \cosh \Delta}, \end{aligned} \quad (43)$$

which, upon regularization, become

$$\begin{aligned} F_+(-i\beta\hbar) &\longrightarrow \frac{1}{2} \int_0^\infty dx e^{-\epsilon x} = \frac{1}{2\epsilon} \\ F_-(-i\beta\hbar) &\longrightarrow \frac{1}{2} \int_0^\infty dx e^{-x(\epsilon + 2 \cosh \Delta)} = \frac{1}{2(\epsilon + 2 \cosh \Delta)}. \end{aligned} \quad (44)$$

This means that

$$Z_{\text{rel}} = \frac{1}{4\epsilon} \pm \frac{1}{4(\epsilon + 2 \cosh \Delta)}. \quad (45)$$

Now, the single particle partition function in this representation is

$$\begin{aligned} Z_1 &= \int d^2r K(r, r; -i\beta\hbar) \Big|_{\mu=m} \\ &\longrightarrow \int_0^\infty dx e^{-\epsilon x} = \frac{1}{\epsilon}, \end{aligned} \quad (46)$$

and the area of the system is regularized to

$$A \longrightarrow \int d^2r \exp\left(-\epsilon \frac{\mu\omega_c r^2}{2\hbar} \operatorname{csch}\Delta\right) = \frac{1}{\epsilon} \lambda_r^2 \left(\frac{\sinh \Delta}{\Delta} \right). \quad (47)$$

Using Eq.(28), one derives the second virial coefficient:

$$\begin{aligned} B_2(\alpha = 0, \Delta, T) &= \frac{A}{2Z_1} [Z_1 - 4Z_{\text{rel}}] \\ &= \mp \frac{1}{4} \lambda_T^2 \left(\frac{\tanh \Delta}{\Delta} \right). \end{aligned} \quad (48)$$

An identical conclusion is reached by a direct expansion of the thermodynamic functions of an ideal quantum gas in a magnetic field. Note that the $\Delta \rightarrow 0$ limit leads to the correct result.

For nonzero α , the situation is more complicated. Defining

$$\begin{aligned} w &= (\epsilon + \cosh \Delta) + \sqrt{(\epsilon + \cosh \Delta)^2 - 1} \\ &= e^\Delta \left(1 + \epsilon \coth \Delta - \frac{1}{2} \epsilon^2 e^{-\Delta} \operatorname{csch}^3 \Delta + \dots \right), \end{aligned} \quad (49)$$

the regularized value of F_\pm takes the form

$$F_\pm(-i\beta\hbar) = \frac{1}{2} e^{-\alpha\Delta} \left[(\epsilon + \cosh \Delta)^2 - 1 \right]^{-1/2} \sum_{n=-\infty}^{\infty} (\pm 1)^n e^{-n\Delta} w^{-|n+\alpha|}. \quad (50)$$

The calculation is carried out by restricting α to one of two intervals: $\alpha \in [-1, 0]$ or $\alpha \in [0, 1]$. The case of arbitrary $\alpha \in \mathbf{R}$ may then be recovered by periodic extension. I shall refer to these two regions as ‘<’ and ‘>’, respectively. The relative coordinate contribution, Z_{rel} , is then

$$Z_{\text{rel}}^< = \frac{w}{w^2 - 1} \left\{ (we^{-\Delta})^\alpha \left(\frac{w^2 e^{-2\Delta}}{w^2 e^{-2\Delta} - 1} \right) + (we^\Delta)^{-\alpha} \left(\frac{1}{w^2 e^{2\Delta} - 1} \right) \right\} \quad (51)$$

and

$$Z_{\text{rel}}^> = \frac{w}{w^2 - 1} \left\{ (we^{-\Delta})^\alpha \left(\frac{1}{w^2 e^{-2\Delta} - 1} \right) + (we^\Delta)^{-\alpha} \left(\frac{w^2 e^{2\Delta}}{w^2 e^{2\Delta} - 1} \right) \right\}. \quad (52)$$

In obtaining B_2 , one must take care to isolate terms of order ϵ as well as terms of order unity before taking the $\epsilon \rightarrow 0$ limit. The second virial coefficient is given by

$$B_2^{>/<}(\alpha, \Delta, T) = \frac{1}{4} \lambda_T^2 \Delta^{-1} \left[\operatorname{ctnh}(\Delta) - 2(\alpha \mp 1) - 2e^{-2(\alpha \mp 1)\Delta} \operatorname{csch} 2\Delta \right], \quad (53)$$

with the top sign applying to region ‘>’, the bottom sign to region ‘<’, and where the $\alpha = 0$ statistics correspond to bosons. One may check explicitly that the special cases $\alpha = 0$ and $\alpha = \pm 1$ lead to the proper ideal quantum gas results of Eq.(48).

Thus, I have derived explicit thermodynamic formulae describing the interpolation of quantum statistics between Bose and Fermi character. The α -dependence

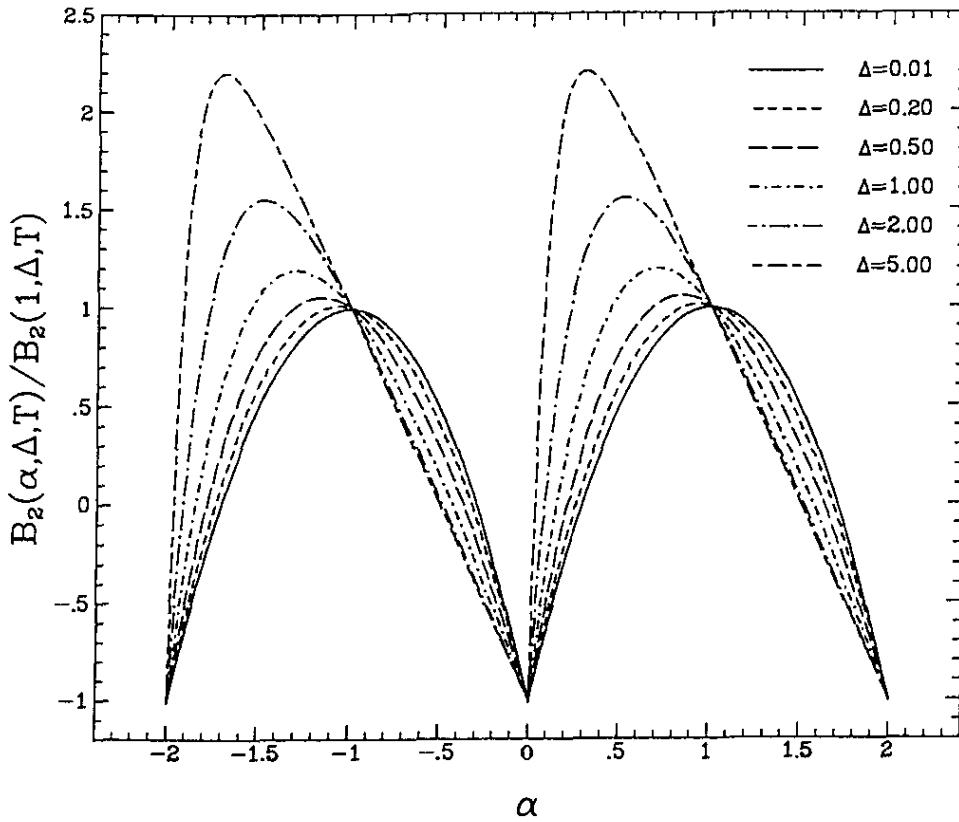


Figure 3. The second virial coefficient $B_2(\alpha, \Delta, T)$ as a function of the statistics determining parameter, α , and scaled by the Fermi value $B_2(1, \Delta, T)$. The asymmetry in the plots is associated with the choice of alignment ($\pm \hat{\mathbf{z}}$) for the magnetic field. All curves are periodic under $\alpha \rightarrow \alpha + 2$, with cusps occurring at even values of α .

of $B_2(\alpha, T)$ at various values of Δ is shown in Fig.[3]. When there is no magnetic field present, Eq.(53) reduces to

$$B_2^{>/<}(\alpha, \Delta, T) \Big|_{\Delta=0} = \frac{1}{4} \lambda_T^2 \left[1 - 2(\alpha \mp 1)^2 \right], \quad (54)$$

which is simply a section of a parabola. The author and coworkers have also obtained the formula of Eq.(54) by a different method¹⁶ which entails performing a Boltzmann-weighted sum of the relative wave function phase shifts $\delta_m(k)$:¹¹

$$B_2(\alpha, T) = \mp \frac{1}{4} \lambda_T^2 - \frac{2}{\pi} \lambda_T^2 \int_0^\infty dk \sum_{m \text{ even/odd}} e^{-\beta \hbar^2 k^2 / 2\mu} \frac{\partial}{\partial k} \delta_m(k), \quad (55)$$

where $\delta_m(k)$ is the phase shift for the m^{th} partial wave.

The virial coefficient $B_2(\alpha, T) \Big|_{\Delta=0}$ is symmetric under the inversion $\alpha \rightarrow -\alpha$. This is clear upon inspection of equations [21] and [24], for the general many-

body Hamiltonian is invariant under the combined operations of inversion in α and complex conjugation. The introduction of a finite magnetic field then breaks this symmetry, which is why the $\Delta \neq 0$ curves of Fig.[3] are skewed toward the left.[†] The $B \neq 0$ Hamiltonian does satisfy

$$H^*(-B, -\alpha) = H(B, \alpha) \quad (56)$$

so all thermodynamic properties will remain invariant under simultaneous inversion of both the field and the statistics determining parameter.

A particularly striking feature is the cusp present at all even integers, those values of α corresponding to bosonic behavior. I will describe the nature of this effect in the case of zero external field by appealing to the relative coordinate energy spectrum of Eq.(20). Recall that if the fiducial statistics are bosonic, only even m states are allowed. Now for small $|\alpha|$, every $m = 0$ partial wave is nondegenerate, as the allowed wavevectors are $k_{l,0} = \Lambda^{-1}x_{|\alpha|,l}$. All states of nonzero angular momentum are, however, *quasidegenerate*: $k_{l,n} \approx k_{l,-n}$. The situation is completely different for α close to one, due to the absence of the odd m states ($m = 1$ in particular), consequently *all* states evidence this quasidegeneracy. This means that whenever α is an odd integer, every state in the sum of Eq.(55) will possess a mate such that their added phase shifts will *cancel* to order (α) , producing the smooth quadratic behavior shown in Fig.[3]. That there are no higher order corrections is an artifact of the ‘free anyon’ model I have considered.

Fractional Statistics and the Quantized Hall Effect

It may seem rather pointless to consider the issue of fractional statistics when our world is clearly not two-dimensional. Even when large anisotropies in couplings effectively reduce the dimensionality of a system, for instance, the condensed matter physicist is ultimately interested in assemblies of electrons and nucleons, for which exotic statistics is a moot issue. The hope, of course, is that certain *elementary excitations* of some quasi-two-dimensional system might behave as if they obeyed fractional statistics. In order that the winding number be a well defined quantity, it seems necessary that any putative anyonic excitations be spatially *localized* and

[†] If the field direction were to be reversed, the curves would then be skewed toward the right.

possess a form factor which is attenuated on some microscopic scale. Thus, extended wave-like objects (*e.g.* phonons), which can be localized only at a high cost in kinetic energy, are poor candidate anyons, and intuition drives one to consider topologically stable local excitations (*e.g.* vortices) and similar entities. Here, I shall discuss the relevance of fractional statistics to the quasiparticle excitations predicted in the fractional quantized Hall effect.[†]

Recall that Laughlin's¹⁷ picture of the fractional quantized Hall effect is based on a discrete set of Jastrow-type *Ansatz* wave functions of the form

$$\Psi_m(z_1, \dots, z_N) = \mathcal{N} \prod_{j < k}^N (z_j - z_k)^m \prod_{l=1}^N e^{-|z_l|^2/4}, \quad (57)$$

where m is an odd integer, $z_j = x_j + iy_j$ is the complex coordinate of particle j , \mathcal{N} is a normalization constant, and where the magnetic length $\ell = \sqrt{\hbar c/eB}$ has been set to unity. Ψ_m describes an incompressible fluid state at filling fraction $\nu = 1/m$ (provided $m \lesssim 70$). The associated quasielectron and quasihole wave functions, respectively, are given by

$$\tilde{\Psi}_m[\boldsymbol{\eta}] = \mathcal{N}(\boldsymbol{\eta}) \prod_{l=1}^N e^{-\bar{z}_l z_l / 4} \prod_{i=1}^N \left(2 \frac{\partial}{\partial z_i} - \bar{\eta} \right) \prod_{j < k}^N (z_j - z_k)^m \quad (58)$$

$$\tilde{\Psi}_m[\boldsymbol{\xi}] = \mathcal{N}(\boldsymbol{\xi}) \prod_{l=1}^N e^{-\bar{z}_l z_l / 4} \prod_{i=1}^N (z_i - \xi) \prod_{j < k}^N (z_j - z_k)^m.$$

The charge of these excitations was also discussed by Laughlin, who employed an argument analogous to that used in deducing the fractional charge of solitons in one-dimensional conductors. He concluded that for $\nu = 1/m$, the quasielectron and quasihole have charges $\mp e^* = \mp e/m$. These excitations are localized within a microscopic region whose size is dictated by the magnetic length and the filling fraction. The excitations described by Eq.(58) are centered at $\mathbf{r} = \boldsymbol{\eta}$ (quasielectron) and $\mathbf{r} = \boldsymbol{\xi}$ (quasiholes), respectively. Roughly speaking, a quasiholes in the incompressible fluid resembles a ‘bubble’ of such a size that $1/m$ of an electron is absent. Less clear, however, is the statistics which the excitations satisfy — Fermi,¹⁷ Bose,¹⁸ and

[†] I shall refer to excitations in the generic sense as ‘quasiparticles’. Specifically, an excitation corresponding to a localized density depletion will be called a ‘quasiholes’, while a local density accumulation will be referred to as a ‘quasielectron’.

fractional¹⁹ statistics have all been proposed. Below, I shall describe a method due to myself and coworkers²⁰ which allows a direct determination of both the charge e^* and statistics α of the excitations.

We determine the charge e^* by evaluating the phase change γ_C accrued by $\tilde{\Psi}_m[\xi]$ as ξ adiabatically moves around a loop of radius R enclosing a flux ϕ . The charge is obtained when γ_C is set equal to the change in phase

$$\gamma_C = \frac{e^*}{\hbar c} \oint \mathbf{A} \cdot d\mathbf{l} = 2\pi \frac{e^* \phi}{e \phi_o}, \quad (59)$$

that an excitation of charge e^* would gain in moving around this loop. The phase γ_C may be calculated via the adiabatic theorem, *treating the quasielectron or quasihole position as an adiabatic parameter.*^{*}

There is a simple yet profound result due to Berry²¹ which relates the adiabatic phase accumulated over such a closed path to a geometric quantity which is *independent* of how slowly the path is traversed. Specifically, let $\{H_\lambda\}$ denote a family of Hamiltonians indexed by the parameter λ . If one executes a closed loop in parameter space sufficiently slowly, the full time-dependent solution to the Schrödinger equation, $\Phi(t)$, will be related to the adiabatic (time-independent) solutions ψ_λ according to[†]

$$\Phi(t) = \exp \left(-i \int^t dt' E(t') \right) e^{i\gamma(t)} \psi(t), \quad (60)$$

where $E(t)$ is the adiabatic energy, $\lambda(t)$ describes the path, and where I have loosely defined the implicit time-dependent functions

$$\begin{aligned} \psi(t) &\equiv \psi_{\lambda(t)} \\ E(t) &\equiv E_{\lambda(t)} \\ \gamma(t) &\equiv \gamma_{\lambda(t)}. \end{aligned} \quad (61)$$

The ‘extra’ phase $\gamma(t)$ satisfies

$$\frac{d}{dt} \gamma(t) = i \langle \psi(t) | \frac{d}{dt} | \psi(t) \rangle. \quad (62)$$

* Although $\tilde{\Psi}_m[\xi]$ is a ‘variational’ wave function, rather than an actual adiabatic wave function, the conclusions obtained are not expected to be sensitive to this inconsistency. We could regard $\tilde{\Psi}_m[\xi]$ to be an exact excited state wave function for a model Hamiltonian.

† I assume that the spectrum of adiabatic wave functions ψ_λ is nondegenerate.

If $\lambda(t)$ describes a *closed* path, then the phase difference

$$\gamma_C \equiv \gamma(T) - \gamma(0) = i \oint dt \langle \psi(t) | \frac{d}{dt} | \psi(t) \rangle \quad (63)$$

has a beautiful geometric interpretation, as noticed by Simon²² (see also Wilczek and Zee²³ and Schiff²⁴).[†]

According to Eq.(58),

$$\frac{d}{dt} \tilde{\Psi}_m[\xi] = \left[\frac{d}{dt} \ln \mathcal{N}[\xi(t)] + \sum_{i=1}^N \frac{d}{dt} \ln |z_i - \xi(t)| \right] \tilde{\Psi}_m[\xi], \quad (64)$$

so that

$$\frac{d}{dt} \gamma(t) = i \frac{d}{dt} \ln \mathcal{N}[\xi] + i \langle \tilde{\Psi}_m[\xi] | \frac{d}{dt} \sum_{i=1}^N \ln(z_i - \xi) | \tilde{\Psi}_m[\xi] \rangle. \quad (65)$$

Using the single particle density in the presence of the quasihole,

$$n_\xi(\mathbf{r}) = \langle \tilde{\Psi}_m[\xi] | \sum_{i=1}^N \delta(\mathbf{r} - \mathbf{r}_i) | \tilde{\Psi}_m[\xi] \rangle, \quad (66)$$

we obtain

$$\frac{d\gamma}{dt} = i \frac{d}{dt} \ln \mathcal{N}[\xi(t)] + i \int d^2 r n_\xi(\mathbf{r}) \frac{d}{dt} \ln |z - \xi(t)|. \quad (67)$$

Since the normalization constant $\mathcal{N}(\xi)$ is a single-valued function of its argument, it will not contribute to the integral expression Eq.(63) for γ_C . We now write $n_\xi(\mathbf{r}) = n_o + \delta n_\xi(\mathbf{r})$, where $n_o = \nu/2\pi$ is the density in the Laughlin ground state Ψ_m , and $\delta n_\xi(\mathbf{r})$ is concentrated about the point $\mathbf{r} = \xi$.^{*} Concerning the n_o term, if ξ is integrated in a clockwise sense around a circle of radius R , one finds

$$\begin{aligned} \oint_{|\xi|=R} dt \frac{d}{dt} \ln |z - \xi(t)| &= \oint_{|\xi|=R} d\xi \frac{1}{\xi - z} \\ &= 2\pi i \theta(R - |z|), \end{aligned} \quad (68)$$

[†] The quantity γ_C is the line integral of a connection form of a $U(1)$ -bundle whose fibers are the adiabatic wave functions. This connection form is given by $\omega = \langle \psi_\lambda | \frac{\partial}{\partial \lambda} | \psi_\lambda \rangle d\lambda$; the fiber metric is defined by the standard inner product. By Stokes' theorem, γ_C can be expressed as a surface integral of the curvature form $\Omega = d\omega + \omega \wedge \omega$. If the parameter space is two-dimensional, then Ω is related to the first Chern form $c_1 = \frac{i}{2\pi} \Omega$, whose integral over the entire parameter space Λ is an integer topological invariant of the bundle.

^{*} The uniform density $n_o = \nu/2\pi$ of the trial ground state and the localized nature of the excitations are easily deduced from the plasma analogy of Laughlin.

where $\theta(x)$ is a step function. Substituting this result into Eq.(67), we find

$$\gamma_C = i \int^R d^2r 2\pi i n_o = -2\pi \langle N \rangle_R = -2\pi\nu\phi/\phi_o, \quad (69)$$

where $\langle N \rangle_R$ denotes the mean number of electrons inside a circle of radius R . We originally expected that corrections to Eq.(69) arising from the $\delta n_\xi(\mathbf{r})$ term would be of order a_0^2/R^2 , where a_0 is the size of the quasi-hole, and that such effects would be irrelevant in the low density limit, where the average separation between quasiholes is much larger than the quasi-hole size. In fact, a much stronger result has been obtained by Haldane²⁵, who has shown that this correction vanishes identically as a consequence of the rotational symmetry of the quasiparticle. Comparing with Eq.(59), we find $e^* = \nu e$, in agreement with Laughlin's result.

A similar analysis shows that the charge of the quasielectron is $e^* = -\nu e$, although one must exercise some caution in dealing with the partial derivative operators in Eq.(58). The adiabatic phase accumulated by a quasielectron is

$$\begin{aligned} \frac{d}{dt} \gamma(t) &= i \frac{d}{dt} \ln \mathcal{N}(\boldsymbol{\eta}) + i \frac{d}{dt} \sum_{i=1}^N \langle \tilde{\Psi}_m[\boldsymbol{\eta}] | \\ &\quad \times \exp \left(-\frac{1}{4} \sum_{l=1}^N |z_l|^2 \right) \ln \left(2 \frac{\partial}{\partial z_i} - \bar{\eta} \right) \exp \left(\frac{1}{4} \sum_{l'=1}^N |z_{l'}|^2 \right) | \tilde{\Psi}_m[\boldsymbol{\eta}] \rangle. \end{aligned} \quad (70)$$

The above matrix element is a many-particle generalization of the Bargmann–Fock space inner product,

$$\langle\langle g | \hat{O} \left(2 \frac{\partial}{\partial z}, z \right) | f \rangle\rangle \equiv \int \frac{dx dy}{2\pi} e^{-\bar{z}z/2} \overline{g(z)} \hat{O} \left(2 \frac{\partial}{\partial z}, z \right) f(z). \quad (71)$$

It is easy to show²⁶ that if the operator \hat{O} is *normal ordered* such that all partial derivatives $\partial/\partial z$ appear to the *left* of all complex coordinates z , then the formal replacement $2\partial/\partial z \rightarrow \bar{z}$ is allowed, *i.e.*

$$\langle\langle g | : \hat{O} \left(2 \frac{\partial}{\partial z}, z \right) : | f \rangle\rangle = \langle\langle g | \hat{O}(\bar{z}, z) | f \rangle\rangle. \quad (72)$$

Making this substitution in Eq.(70), one recovers the result of Eq.(67), with $z \rightarrow \bar{z}$ and $\xi \rightarrow \bar{\eta}$. The charge of the quasielectron follows immediately.

To determine the statistics of the excitations, we consider the state with quasi-holes at ξ and ξ' ,

$$\tilde{\Psi}_m[\xi, \xi'] = \mathcal{N}(\xi, \xi') \prod_{i=1}^N [(z_i - \xi)(z_i - \xi')] \Psi_m. \quad (73)$$

As above, we adiabatically carry ξ around a closed loop of radius R . If ξ' lies outside the circle $|\xi| = R$ by a distance $d \gg a_0$, the above analysis is unchanged, i.e. $\gamma_C = -2\pi\nu\phi/\phi_0$. If, however, ξ' lies inside this loop and $R - |\xi'| \ll a_0$, there is a deficit in $\langle N \rangle_R$ of $-\nu$, and the phase accrued is $\gamma'_C = \gamma_C + 2\pi\nu$. Therefore, when a quasi-hole adiabatically encircles another quasi-hole, an extra ‘statistical phase’

$$\Delta\gamma_C = 2\pi\nu \quad (74)$$

is accumulated. Again, an analogous result holds for quasielectron. For the case of the filled Landau level, $\nu = 1$, $\Delta\gamma_C = 2\pi$, and the phase obtained upon interchanging quasiholes is $\Delta\gamma_C/2 = \pi$, corresponding to Fermi statistics. For ν noninteger, we identify the quantity $\Delta\gamma_C/2\pi$ with the statistics determining parameter, α , and the excitations obey fractional statistics, in agreement with the conclusion of Halperin.¹⁹ Clearly, when ν is nonintegral, the change of phase $\Delta\gamma_C$ accumulated when a *third* particle is in the vicinity will depend on the adiabatic path taken by the excitations as they are interchanged, and the permutation definition used for Bose and Fermi statistics no longer suffices.

It has been noted by Su²⁷ and others that the above derivation is sensitive only to certain very general properties of Ψ_m and $\tilde{\Psi}_m[\xi]$, and that the same conclusions apply whenever the ground state Ψ is of uniform density and the quasiparticles are localized and of the form given in Eq.(58).

There is an unfortunate confusion involving the nature of the elementary excitations which has found its way into the literature.²⁸ *Although they may obey fractional statistics, the quasiparticles do not carry a magnetic flux $\phi = \nu\phi_0$ as do Wilczek’s charged particle-flux tube composites.* The exotic statistics derive from the sharp fractional charge of the excitations and the incompressibility of the ground state, and there are *no* trapped flux lines in the Hall plasma.

While the derivation is certainly suggestive, its interpretation is not completely clear. For instance, what is the meaning of treating the quasiparticle location as an adiabatic parameter? I believe that this *is* the correct procedure, and I will now present a simple analogous result which seems to justify this assumption. Consider the single particle coherent state $|\mathbf{R}\rangle$,

$$\varphi_{\mathbf{R}}(\mathbf{r}) = \langle \mathbf{r} | \mathbf{R} \rangle = \frac{1}{\sqrt{2\pi}} e^{-(\mathbf{r}-\mathbf{R})^2/4} e^{-i\mathbf{r} \times \mathbf{R} \cdot \hat{\mathbf{z}}/2}. \quad (75)$$

If the guiding center is adiabatically transported around a closed loop, it is easy to show that the Berry phase satisfies

$$\frac{d}{dt} \gamma(t) = i \cdot \mathbf{R}(t) + \frac{d}{dt} |\mathbf{R}(t)| = -\frac{1}{\phi_0} \frac{d\phi}{dt}, \quad (76)$$

and adiabatically dragging the electron about a closed path leads to a cumulative Berry phase of $-\phi/\phi_0$, from which the correct electron charge is recovered. Thus, in this trivial example, treating the position as an adiabatic parameter leads to the appropriate results.

Halperin's 'Pseudo Wave Function'

Halperin¹⁹ originally used fractional statistics to derive the hierarchy scheme of rational filling fractions based on successive levels of Laughlin condensates. In this manner, states of density $\nu \neq 1/m$ could be described as an underlying primitive ($1/m$) Laughlin condensate plus a condensate of quasielectrons or quasiholes. If such a composite state does not fully saturate the density, quasiparticles of the 'second level' can be added and themselves condense, and so on, until the desired rational fraction $\nu = p/q$ is reached. As expected, the gap decreases sharply as one climbs up this tree of hierarchical composite states, and experimental conditions will determine how far up the tree the system may progress and still remain condensed.

Each excitation condensate at a level s of the hierarchy is defined by a 'pseudo wave function', Φ_s , whose arguments are the positions $Z_j = X_j \pm iY_j$ of the excitations (the plus sign being taken for quasiholes, and the minus sign for quasielectrons, at every level except $s = 0$). I shall henceforth refer to the condensed entities at stage s in the hierarchy as ' s -particles' and thereby distinguish them from ' s -excitations', which will be used to refer to defects in the s -condensate. The explicit form of Φ_s is

$$\begin{aligned} \Phi_s(\mathbf{R}_1, \dots, \mathbf{R}_N) &= P_s(Z_1, \dots, Z_N) Q_s(Z_1, \dots, Z_N) \prod_{l=1}^N e^{-|q_s|Z_l Z_l/4} \\ P_s(Z_1, \dots, Z_N) &= \prod_{j < k}^N (Z_j - Z_k)^{2p_{s+1}} \\ Q_s(Z_1, \dots, Z_N) &= \prod_{j < k}^N (Z_j - Z_k)^{-\sigma_{s+1}/m_s}, \end{aligned} \quad (77)$$

where the excitation charge is $e_s = \pm q_s e$, σ_{s+1} denotes the type of s -particle (+1 for electron-like objects, and -1 for hole-like objects), and $m_s > 1$ is a rational number which will be determined by iteration in s . The polynomial Q_s determines the statistics of the entities described by the pseudo wave function — interchange of two s -particles in Eq.(77) leads to a phase of $e^{\pm i\pi m_s}$. With p_{s+1} restricted to the positive integers, multiplication by the (symmetric) polynomial P_s will introduce additional correlations in Φ_s without altering the statistics. For $m_s \neq 1$, the function Φ_s is multiple-valued and thus is formally well defined only on the universal covering of the configuration space for the M s -particles. This last feature of Φ_s is not a particularly important one, however, as far as Halperin's analysis goes.

Halperin interprets Φ_s as the probability amplitude for finding an $(s-1)$ -excitation at each of the N positions $(\mathbf{R}_1, \dots, \mathbf{R}_N)$, provided that these $(s-1)$ -excitations are sufficiently well separated so that their cores do not overlap.[†] Aside from this latter restriction, Φ_s may be considered to be a conventional, quantum mechanical wave function.

Invoking Laughlin's plasma analogy, the number density of s -particles in the state Φ_s is found to be

$$\begin{aligned} n_s &= |q_s|/2\pi m_{s+1} \\ m_{s+1} &\equiv 2p_{s+1} - \sigma_{s+1}/m_s. \end{aligned} \tag{78}$$

The actual charge of the s -particles at level s is $\tilde{e}_s = \sigma_{s+1} q_s e$. Clearly the true electronic filling fraction, ν_{s+1} , in this state is

$$\nu_{s+1} = \nu_s + \sigma_{s+1} q_s |q_s|/m_{s+1}. \tag{79}$$

By applying the quasielectron (quasihole) operators of Eq.(58),

$$\begin{aligned} \hat{O}_+[\bar{\eta}] &= \prod_{i=1}^N \left(2 \frac{\partial}{\partial Z_i} - \bar{\eta} \right) \\ \hat{O}_-[\xi] &= \prod_{i=1}^N (Z_i - \xi), \end{aligned} \tag{80}$$

one accumulates a surplus (deficit) of $1/m_{s+1}$ s -particles about the point $Z = \eta$ ($Z = \xi$). These s -excitations then comprise the fundamental $(s+1)$ -particles of

[†] Since the s -particles are in general not point particles, as are the electrons of the primitive Laughlin condensates, the pseudo wave function has no meaning when any two s -particles lie too close to each other.

the next level -- this procedure defines the basic link in the hierarchy chain. This reasoning leads to an iterative equation for the q_s :

$$q_{s+1} = \sigma_{s+1} q_s / m_{s+1}. \quad (81)$$

At the base of the hierarchy, the initial conditions are $\nu_0 = 0$, $q_0 = m_0 = \sigma_1 = 1$. By specifying a sequence $\{p_s, \sigma_s\}$, the formulae [78-81] then determine a corresponding sequence of rational filling fractions. By expressing the filling fractions in terms of continued fractions, one recovers the hierarchy scheme of Haldane.¹⁸

Halperin has derived the starting expression of Eq.(77) in a somewhat systematic fashion by introducing additional correlations between *pairs* of electrons in Laughlin's primitive trial state. However, it may be more than merely a curious fact that the normalization integral, $\mathcal{N}(\xi_1, \dots, \xi_M)$, of the excited state pseudo wave function

$$\tilde{\Phi}_s(\xi_1, \dots, \xi_M; \mathbf{R}_1, \dots, \mathbf{R}_N) = \mathcal{N}_s(\xi_1, \dots, \xi_M) \prod_{\alpha=1}^M \hat{O}_-[\xi_\alpha] \Phi_s(\mathbf{R}_1, \dots, \mathbf{R}_N) \quad (82)$$

is actually *equal*, apart from an overall constant, to the modulus of the next pseudo wave function in the hierarchy, $\Phi_{s+1}|_{p_{s+2}=0}$.

The normalization integral for quasiholes is given by

$$|\mathcal{N}_s(\xi_1, \dots, \xi_M)|^2 = \int \prod_{i=1}^N d^2 R_i \prod_{i=1}^N \prod_{\alpha=1}^M |Z_i - \xi_\alpha|^2 |\Phi_s(\mathbf{R}_1, \dots, \mathbf{R}_N)|^2. \quad (83)$$

Differentiation with respect to ξ_β yields

$$\begin{aligned} -2 \frac{\partial}{\partial \xi_\beta} \ln |\mathcal{N}_s| &= \int \prod_{i=1}^N d^2 R_i \left(\sum_{j=1}^N \frac{1}{\xi_\beta - Z_j} \right) |\tilde{\Phi}_s|^2 \\ &= - \int d^2 r \left(\frac{1}{Z - \xi_\beta} \right) n(\mathbf{R}), \end{aligned} \quad (84)$$

where $n(\mathbf{R})$ is the number density of the fundamental level s particles. If the excitations are few in number, this equation may be analyzed using the same reasoning employed in calculating Berry's phase for the Laughlin quasihole. The result is

$$\begin{aligned} -2 \frac{\partial}{\partial \xi_\beta} \ln |\mathcal{N}_s| &= \pi n_s \bar{\xi}_\beta + \frac{1}{2m_{s+1}} \sum_{\substack{\alpha=1 \\ \alpha \neq \beta}}^M \left(\frac{1}{\xi_\alpha - \xi_\beta} \right) \\ |\mathcal{N}_s(\xi_1, \dots, \xi_M)| &= \mathcal{A} \prod_{\alpha < \beta}^M |\xi_\alpha - \xi_\beta|^{1/m_{s+1}} \prod_{\gamma=1}^M \exp(-|q_s| \bar{\xi}_\gamma \xi_\gamma / 4m_{s+1}), \end{aligned} \quad (85)$$

where \mathcal{A} is a constant. Had I considered quasielectron excitations, the statistics determining exponent, $1/m_{s+1}$, would have appeared with the opposite sign — the type of excitation thus defines the quantity σ_{s+2} . Thus, the normalization integral at level s in the hierarchy is simply the modulus of the pseudo wave function at level $s + 1$, with $p_{s+2} = 0$. Again, one may introduce additional correlations without altering the statistics by choosing a nonzero integer for p_{s+2} , thereby recovering the general expression for the level $s + 1$ pseudo wave function.

Laughlin was kind enough to tell me that this result is trivial.²⁹ Consider, for example, the expression

$$\Upsilon(\xi_1, \dots, \xi_M; r_1, \dots, r_N) \equiv \prod_{\alpha < \beta}^M |\xi_\alpha - \xi_\beta|^{2/m} \prod_{j < k}^N |z_j - z_k|^{2m} \\ \times \prod_{i=1}^N \prod_{\delta=1}^M |z_i - \xi_\delta|^2 \prod_{\gamma=1}^M e^{-\bar{\xi}_\gamma \xi_\gamma / 2m} \prod_{l=1}^N e^{-\bar{z}_l z_l / 2}, \quad (86)$$

which may be interpreted à la Laughlin as the classical probability distribution of a plasma consisting of N particles of charge m and M particles of charge 1 at a temperature $T = m$ interacting via logarithmic potentials, including the usual neutralizing background. If one integrates out the ‘electrons’ (charge m particles) in the dilute ξ -particle limit, one expects to obtain the partition function for a constrained system (constrained in that the ξ coordinates are fixed), and this should resemble $e^{-M\mu_\xi/T}$, where μ_ξ is the chemical potential for the ξ -particles. But then

$$\int \Upsilon d^2r_1 \cdots d^2r_N = |\mathcal{N}(\xi_1, \dots, \xi_M)|^{-2} \prod_{\alpha < \beta}^M |\xi_\alpha - \xi_\beta|^{2/m} \prod_{\gamma=1}^M e^{-\bar{\xi}_\gamma \xi_\gamma / 2m} \\ = e^{-M\mu_\xi/T}, \quad (87)$$

from which the pseudo wave function modulus is obtained immediately. This result holds in the dilute gas limit where the ξ -particles are equally likely to be anywhere in the plasma except within a Debye length of each other. At any rate, the result is the same as that of Eq.(85).

This manner in which the pseudo wave function is derived is very suggestive of something, although what that something is I do not know.

Selection Rules for Anyon Production

Tao³⁰ has used the fractional statistics of the excitations to argue for the odd

denominator rule. This may be putting the cart a bit before the proverbial horse, however Tao's analysis does lead to some selection rules for anyon production. Suppose, for example, that an electron decays into a number of anyons, or, more generally, that p electrons decay into q identical anyons. Without loss of generality, p and q can be taken to be relatively prime. Consider now two groups of q such anyons apiece. Interchange of these groups must give a phase

$$\Delta\gamma = (-1)^{p^2} = e^{i\pi p^2}. \quad (88)$$

because each group derives from p electrons. On the other hand, if the two assemblies are viewed as consisting of anyons, the statistical phase will be

$$\Delta\gamma = e^{i\alpha\pi q^2} = e^{i\pi pq}, \quad (89)$$

if the statistics determining parameter is tied to the charge by $\alpha = e^*/e$. The two formulae of Eqs.(88,89) are incompatible if q is even, leading to the selection rule that fermion-to-anyon decays must proceed via an odd- p , odd- q channel or an even- p , odd- q channel. If the parent particle is a boson, one finds that the boson-to-anyon decays may proceed by either an odd- p , even- q channel or an even- p , odd- q channel.

While the virial coefficient calculations described earlier are somewhat relevant to the fractional quantum Hall effect, they do not take into account the Coulomb interaction between the excitations — I have also neglected effects due to the finite size of the quasielectron and quasihole. In fact, in real systems, the one-body potential along the semiconductor interface is presumed strong enough to pin the excitations (this is the qualitative picture of the Hall plateaux). It is conceivable that a charged impurity may bind m fractionally charged excitations ($e^* = \pm e/m$) and thus form an ‘atom’. The energy levels, or shell structure, of this atom would presumably be sensitive to the statistics of the excitations.

One might also look for fractional statistics numerically. For example, if one examines FQHE clusters away from $\nu = \frac{1}{3}$, the ground state will resemble a quiescent Laughlin-type liquid with a certain number of quasiparticle excitations present. The energy levels could then be compared to those of a system of noninteracting quasiparticles, for which the statistics of the excitations will have definite consequences.

Summary

In two dimensions, a peculiar formulation of quantum statistics can be implemented if one associates to each particle a fictitious charge \tilde{e} and a flux tube of strength $\phi = \alpha hc/2\tilde{e}$. Evaluating the many-body propagator by a path integral, a fictitious Aharonov-Bohm phase is accrued when particles wind around one another; physical quantities are found to be periodic under $\alpha \rightarrow \alpha + 2$. This procedure defines one among a continuum of possible $(2+1)$ -dimensional quantum theories indexed by α , which, due to the aforementioned periodicity, may be restricted to the range $[0, 2)$. Conventional statistics arise when α is an integer: $\alpha = 0$ for bosons, and $\alpha = 1$ for fermions. For general α , the particles are referred to as *anyons*.

Due to the long-ranged nature of this ‘statistical’ interaction, it is in general difficult to calculate anything interesting within this formalism. In this chapter, I have discussed the solution to the two-anyon problem, which may be used to understand the lowest order quantum corrections to anyon thermodynamics. Specifically, the second virial coefficient, $B_2(\alpha, T)$ was evaluated, and was found to evidence nonanalytic behavior as a function of the statistics determining parameter, nonetheless properly interpolating between its corresponding Bose and Fermi values.

I have also elucidated some of the arguments as to why the Laughlin quasiparticles in the fractional quantized Hall effect should obey fractional statistics. In this case, the statistics determining parameter is itself determined by the ground state filling fraction: $\alpha = \nu \equiv nhc/eB$. While such a conclusion is in conformity with the hierarchical condensate scheme,^{18,19} no simple experimental (or numerical) test has yet been performed which could test this prediction.

Acknowledgements

I have greatly benefited from discussions with R. MacKenzie, J. R. Schrieffer, and F. Wilczek. This work is excerpted from my Ph.D. thesis (University of California at Santa Barbara, 1986) I am grateful to AT&T Bell Laboratories for their support while this work was in progress.

APPENDIX A: PROPAGATOR IN THE PRESENCE OF A FLUX TUBE

In this appendix, I shall discuss the method of path integration in polar coordinates and its application to the quantum mechanical propagator of a free particle in the presence of a flux tube.

What is desired is an expression for the propagator

$$K^\alpha(\mathbf{r}'', \mathbf{r}'; T) = \int \mathcal{D}\mathbf{r}(t) \exp [iS(\mathbf{r}'', \mathbf{r}')/\hbar], \quad (\text{A.1})$$

with $T \equiv t'' - t'$; the label α will hereafter be suppressed. The technique involves isolating the contribution of a particular homotopy sector to the path integral (here, the homotopy index is simply the winding number of a path about the origin). The constraint

$$\int_{t'}^{t''} dt \dot{\theta} = \phi \quad (\text{A.2})$$

is imposed on the propagator by weighing each path with a δ -function:

$$K_\phi(\mathbf{r}'', \mathbf{r}'; T) = \int \mathcal{D}\mathbf{r}(t) e^{iS(\mathbf{r}'', \mathbf{r}')/\hbar}. \quad (\text{A.3})$$

The above expression will be nonzero *only* for those ϕ which satisfy $\phi = \theta'' - \theta' + 2\pi n$, and the complete propagator is clearly

$$K(\mathbf{r}'', \mathbf{r}'; T) = \int d\phi K_\phi(\mathbf{r}'', \mathbf{r}'; T). \quad (\text{A.3})$$

The quantity K_ϕ is known as the *constrained propagator*. Using an integral representation of the δ -function, one finds

$$K_\phi(\mathbf{r}'', \mathbf{r}'; T) = \int_{-\infty}^{\infty} \frac{d\lambda}{2\pi} e^{i\lambda\phi} \int \mathcal{D}\mathbf{r}(t) e^{iS_\lambda(\mathbf{r}'', \mathbf{r}')/\hbar}$$

$$S_\lambda(\mathbf{r}'', \mathbf{r}') = \int_{t'}^{t''} dt [L(\mathbf{r}, \dot{\mathbf{r}}; t) - \lambda \hbar \dot{\theta}] \equiv \int_{t'}^{t''} dt L_\lambda(\mathbf{r}, \dot{\mathbf{r}}, t). \quad (\text{A.4})$$

The relative coordinate problem in which I am interested is defined by the Lagrangian

$$L(\mathbf{r}, \dot{\mathbf{r}}) = \frac{1}{2}\mu\dot{\mathbf{r}}^2 - \frac{1}{2}\omega_c^2 r^2 \dot{\theta} - \alpha\dot{\theta}. \quad (\text{A.5})$$

This is cast into a more manageable form by redefining θ :

$$\vartheta \equiv \theta - \frac{1}{2}\omega_c t$$

$$L_\lambda = \frac{1}{2}\mu\dot{\mathbf{r}}^2 - \frac{1}{8}\mu r^2 \omega_c^2 - (\alpha + \lambda)\dot{\vartheta} - \frac{1}{2}(\alpha + \lambda)\hbar\omega_c. \quad (\text{A.6})$$

The path integral is calculated by the usual discretization procedure:

$$\begin{aligned}
 K_\phi(\mathbf{r}'', \mathbf{r}'; T) &= \lim_{N \rightarrow \infty} \left(\frac{\mu}{2\pi i \hbar \epsilon} \right)^N \int \frac{d\lambda}{2\pi} e^{i\lambda \phi} \\
 &\quad \times \int d^2 r_1 \cdots \int d^2 r_N \exp \left[i \sum_{j=1}^{N+1} S_\lambda(\mathbf{r}_j, \mathbf{r}_{j-1}) / \hbar \right], \\
 S_\lambda(\mathbf{r}_j, \mathbf{r}_{j-1}) &\approx \epsilon L_\lambda \left[\frac{1}{\epsilon} (\mathbf{r}_j - \mathbf{r}_{j-1}), \mathbf{r}_j \right] \\
 &= \frac{\mu}{2\epsilon} \left(r_j^2 + r_{j-1}^2 - \frac{1}{4} r_j^2 \omega_c^2 \epsilon^2 \right) \\
 &\quad - \frac{\mu}{\epsilon} r_j r_{j-1} \cos(\vartheta_j - \vartheta_{j-1}) - \bar{\lambda} \hbar (\vartheta_j - \vartheta_{j-1}) - \frac{1}{2} \bar{\lambda} \hbar \omega_c \epsilon,
 \end{aligned} \tag{A.7}$$

with $\bar{\lambda} = \lambda + \alpha$, $\epsilon = T/(N+1)$, $\mathbf{r}_0 = \mathbf{r}'$, and $\mathbf{r}_{N+1} = \mathbf{r}''$. What makes this problem difficult is the periodicity of the cosine and the finite limits on the ϑ integrals. To properly account for all terms of order ϵ^2 , the approximation

$$\cos(\Delta\vartheta) - a\epsilon \Delta\vartheta \approx \cos(\Delta\vartheta - a\epsilon) + \frac{1}{2} a^2 \epsilon^2 \tag{A.8}$$

is employed. This yields

$$\begin{aligned}
 iS(\mathbf{r}_j, \mathbf{r}_{j-1}) / \hbar &= \frac{i\mu}{2\hbar\epsilon} (r_j^2 + r_{j-1}^2 - \frac{1}{4} r_j^2 \omega_c^2 \epsilon^2) - \frac{1}{2} \bar{\lambda} \hbar \omega_c \epsilon - \frac{1}{2} i \bar{\lambda}^2 \frac{\hbar\epsilon}{\mu r_j r_{j-1}} \\
 &\quad - \frac{i\mu r_j r_{j-1}}{\hbar\epsilon} \cos \left(\vartheta_j - \vartheta_{j-1} - \bar{\lambda} \frac{\hbar\epsilon}{\mu r_j r_{j-1}} \right).
 \end{aligned} \tag{A.9}$$

Further progress is made by appealing to the generating function

$$\exp \left[\frac{1}{2} z \left(s + s^{-1} \right) \right] = \sum_{m=-\infty}^{\infty} s^m I_m(z) \tag{A.10}$$

for the modified Bessel function of the first kind. $I_\nu(x)$ behaves asymptotically as

$$I_\nu(x) \sim \frac{e^x}{\sqrt{2\pi x}} \left\{ 1 - \frac{1}{2} (\nu^2 - \frac{1}{4}) \frac{1}{x} + \dots \right\}, \tag{A.11}$$

which leads to the following limit:

$$\lim_{\epsilon \rightarrow 0} \exp \left(-\frac{1}{2} i \bar{\lambda}^2 \frac{\epsilon}{x} \right) \exp \left[-i \frac{x}{\epsilon} \cos \left(\theta - \bar{\lambda} \frac{\epsilon}{x} \right) \right] = \sum_{m=-\infty}^{\infty} e^{im\theta} I_{|m+\bar{\lambda}|} \left(\frac{x}{i\epsilon} \right). \tag{A.12}$$

The constrained propagator, then, resembles

$$\begin{aligned}
 K_\phi(\mathbf{r}'', \mathbf{r}'; T) &= \lim_{N \rightarrow \infty} \left(\frac{\mu}{2\pi i \hbar \epsilon} \right)^N e^{-i\alpha\phi} \int \frac{d\bar{\lambda}}{2\pi} e^{i\bar{\lambda}(\phi - \frac{1}{2} \omega_c T)} \\
 &\quad \times \int d^2 r_1 \cdots \int d^2 r_N \prod_{j=1}^{N+1} \sum_{m_j=-\infty}^{\infty} \left\{ e^{im_j(\vartheta_j - \vartheta_{j-1})} \right. \\
 &\quad \times \left. \exp \left[\frac{i\mu}{2\hbar\epsilon} (r_j^2 + r_{j-1}^2 - \frac{1}{4} r_j^2 \omega_c^2 \epsilon^2) \right] I_{|m_j + \bar{\lambda}|} \left(\frac{\mu r_j r_{j-1}}{i\hbar\epsilon} \right) \right\}.
 \end{aligned} \tag{A.13}$$

At this point, the ϑ_j integrations may be carried out. The expression partly collapses due to the orthonormality relation

$$\int_{-\pi}^{\pi} \frac{d\vartheta}{2\pi} e^{i(m'-m)\vartheta} = \delta_{m,m'}, \quad (\text{A.14})$$

and one obtains

$$\begin{aligned} K_\phi(\mathbf{r}'', \mathbf{r}'; T) &= \lim_{N \rightarrow \infty} \left(\frac{\mu}{i\hbar\epsilon} \right)^N e^{-i\alpha\phi} \int \frac{d\lambda}{2\pi} e^{i\lambda(\phi - \frac{1}{2}\omega_c T)} \\ &\times \int d^2 r_1 \dots \int d^2 r_N \prod_{j=1}^{N+1} \left\{ \sum_{m=-\infty}^{\infty} e^{im(\vartheta_{N+1} - \vartheta_0)} \right. \\ &\times \left. \exp \left[\frac{i\mu}{2\hbar\epsilon} \left(r_j^2 + r_{j-1}^2 - \frac{1}{4} r_j^2 \omega_c^2 \epsilon^2 \right) \right] I_{|m+\lambda|} \left(\frac{\mu r_j r_{j-1}}{i\hbar\epsilon} \right) \right\}. \end{aligned} \quad (\text{A.15})$$

This still looks sufficiently ugly so as to be rather intimidating. However, a truly remarkable result due to Peak and Inomata¹³ saves the day. By iterating the formula

$$\int_0^\infty dx e^{i\sigma x} I_\nu(-ia\sqrt{x}) I_\nu(-ib\sqrt{x}) = \frac{i}{\sigma} \exp \left[\frac{-i(a^2 + b^2)}{4\sigma} \right] I_\nu \left(\frac{ab}{2\sigma} \right) \quad (\text{A.16})$$

N times ($\text{Re}(\nu) > -1$), one finds

$$\begin{aligned} \int_0^\infty \prod_{k=1}^N dr_k r_k \exp \left(ia \sum_{j=1}^N r_j^2 \right) \prod_{j=1}^{N+1} I_\nu(-ib r_j r_{j-1}) = \\ \prod_{k=1}^N \left(\frac{i}{2a_k} \right) \exp \left\{ -i \left[r'^2 \sum_{j=1}^N \frac{b_j^2}{4a_j} + r''^2 \frac{b^2}{a_{N+1}} \right] \right\} I_\nu(-ib_{N+1} r' r''), \end{aligned} \quad (\text{A.17})$$

with

$$\begin{aligned} a_1 &= a, & a_{j+1} &= a - \frac{b^2}{4a_j}, & j &\geq 1, \\ b_1 &= b, & b_{j+1} &= \prod_{k=1}^j \frac{b}{2a_k}, & j &\geq 1. \end{aligned} \quad (\text{A.18})$$

Peak and Inomata showed that the $N \rightarrow \infty$ limit leads to

$$\begin{aligned} \lim_{N \rightarrow \infty} \left(\prod_{j=1}^N \frac{b}{2a_j} \right) &= \frac{1}{2} \mu \omega_c \csc(\frac{1}{2}\omega_c T) \\ \lim_{N \rightarrow \infty} \left(\frac{1}{2}b - \sum_{j=1}^N \frac{b_j^2}{4a_j} \right) &= \frac{1}{4} \mu \omega_c \operatorname{ctn}(\frac{1}{2}\omega_c T) \\ \lim_{N \rightarrow \infty} \left(\frac{1}{2}b - \frac{b^2}{4a_N} \right) &= \frac{1}{4} \mu \omega_c \operatorname{ctn}(\frac{1}{2}\omega_c T) \end{aligned} \quad (\text{A.19})$$

where

$$\begin{aligned} a &\equiv \frac{\mu}{\epsilon} \left(1 - \frac{1}{8}\omega_c^2\epsilon^2\right) \\ b &= \frac{\mu}{\epsilon}. \end{aligned} \quad (\text{A.20})$$

Using this result, the constrained propagator may finally be written as

$$\begin{aligned} K_\phi(\mathbf{r}'', \mathbf{r}'; T) &= e^{-i\alpha\phi} \int \frac{d\lambda}{2\pi} e^{i\lambda(\phi - \frac{1}{2}\omega_c T)} \sum_{m=-\infty}^{\infty} e^{im(\theta'' - \theta' - \frac{1}{2}\omega_c T)} Q_{|m+\lambda|}(\mathbf{r}'', \mathbf{r}'; T) \\ &= e^{-i\alpha\phi} \sum_{n=-\infty}^{\infty} \delta(\theta'' - \theta' - \phi + 2\pi n) \int d\bar{\lambda} e^{i\bar{\lambda}(\phi - \frac{1}{2}\omega_c T)} Q_{|\bar{\lambda}|}(\mathbf{r}'', \mathbf{r}'; T) \end{aligned} \quad (\text{A.21})$$

where the radial function $Q_\beta(r'', r'; T)$ is defined to be

$$\begin{aligned} Q_\beta(r'', r'; T) &\equiv \frac{\mu\omega_c}{4\pi i\hbar} \csc\left(\frac{1}{2}\omega_c T\right) \\ &\times \exp\left[\frac{i\mu\omega_c}{4\hbar} \operatorname{ctn}\left(\frac{1}{2}\omega_c T\right)(r'^2 + r''^2)\right] \\ &\times I_\beta\left(\frac{\mu\omega_c}{2i\hbar} r' r'' \csc\left(\frac{1}{2}\omega_c T\right)\right). \end{aligned} \quad (\text{A.22})$$

In the limit of zero field, the radial function becomes

$$\lim_{\omega_c \rightarrow 0} Q_\beta(r'', r'; T) = \frac{\mu}{2\pi i\hbar T} \exp\left[\frac{i\mu}{2\hbar T}(r'^2 + r''^2)\right] I_\beta\left(\frac{\mu r' r''}{i\hbar T}\right). \quad (\text{A.23})$$

As advertised, the winding constraint $\int dt \dot{\theta} = \phi$ renders the expression (A.21) for K_ϕ zero unless $\phi = \theta'' - \theta' + 2\pi n$. The complete propagator, K , is obtained by integrating over this constraint:

$$K(\mathbf{r}'', \mathbf{r}'; T) = e^{-i\alpha\omega_c T/2} \sum_{n=-\infty}^{\infty} e^{in(\theta'' - \theta' - \frac{1}{2}\omega_c T)} Q_{|n+\alpha|}(\mathbf{r}'', \mathbf{r}'; T). \quad (\text{A.24})$$

Of course, taking both $\omega_c \rightarrow 0$ and $\alpha = 0$ in Eq.(A.24) recovers the familiar free particle propagator,

$$\lim_{\omega_c \rightarrow 0} K(\mathbf{r}'', \mathbf{r}'; T) \Big|_{\alpha=0} = \frac{\mu}{2\pi i\hbar T} \exp\left[i(\mathbf{r}'' - \mathbf{r}')^2/2\hbar T\right]. \quad (\text{A.25})$$

I wish to stress that this appendix is meant as a technical review and that all the formulae herein have been derived elsewhere.

REFERENCES

- [1] Frank Wilczek, *Phys. Rev. Lett.* **48**, 1144 (1982).
- [2] Yong-Shi Wu, *Phys. Rev. Lett.* **53**, 111 (1984).
- [3] Michael G. G. Laidlaw and Cécile Morette DeWitt, *Phys. Rev. D* **3**, 6 (1971).
- [4] Frank Wilczek, *Phys. Rev. Lett.* **49**, 957 (1982).
- [5] Yong-Shi Wu, *Phys. Rev. Lett.* **52**, 2103 (1984).
- [6] L. Schulman, *Techniques and Applications of Path Integration*, (Wiley, New York, 1981). Chapter 27 contains material on coherent state path integration, while chapter 23 gives an excellent account of path integration on multiply connected spaces.
- [7] Edward Fadell and James Van Buskirk, *Duke Math J.* **29**, 243 (1962).
- [8] Joan S. Birman, *Comm. Pure and App. Math.* **22**, 41 (1969).
- [9] See, for example, R. Fox and L. Neuwirth, *Math. Scand.* **10**, 119 (1962).
- [10] D. C. Hickernell, E. O. McLean, and O. E. Vilches, *Phys. Rev. Lett.* **28**, 789 (1972); M. Bretz, J. G. Dash, D. C. Hickernell, E. O. McLean, and O. E. Vilches, *Phys. Rev. A* **8**, 1589 (1973) and *Phys. Rev. A* **9**, 2814 (1974).
- [11] J. G. Dash, *Films on Solid Surfaces*, (Academic Press, New York, 1975).
- [12] S. F. Edwards and Y. V. Gulyaev, *Proc. Roy. Soc. London A* **279**, 229 (1964).
- [13] D. Peak and A. Inomata, *J. Math. Phys.* **10**, 1422 (1969).
- [14] Akira Inomata and Vijay A. Singh, *J. Math. Phys.* **19**, 2318 (1978); Christopher C. Gerry and Vijay A. Singh, *Phys. Rev. D* **20**, 2550 (1979); C. C. Gerry and V. A. Singh, *Il Nuovo Cimento* **73B**, 161 (1983).
- [15] I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products*, (Academic Press, New York, 1980). The result (6.611.4) is particularly useful for those interested in the material contained in appendix 5A of this thesis.
- [16] Daniel P. Arovas, Robert Schrieffer, Frank Wilczek, and A. Zee, *Nucl. Phys.* **B251 [FS13]**, 117 (1985).

- [17] R. B. Laughlin, *Phys. Rev. Lett.* **50**, 1395 (1983).
- [18] F. D. M. Haldane, *Phys. Rev. Lett.* **51**, 605 (1983).
- [19] B. I. Halperin, *Phys. Rev. Lett.* **52**, 1583 (1984).
- [20] Daniel Arovas, J. R. Schrieffer, and Frank Wilczek, *Phys. Rev. Lett.* **53**, 722, 1984.
- [21] M. V. Berry, *Proc. Roy. Soc. London* **A392**, 45 (1984).
- [22] B. Simon, *Phys. Rev. Lett.* **51**, 2167 (1983).
- [23] Frank Wilczek and A. Zee, *Phys. Rev. Lett.* **52**, 2111 (1984).
- [24] L. Schiff, *Quantum Mechanics*, (McGraw-Hill, New York, 1955), p. 290.
- [25] F. D. M. Haldane, private communications.
- [26] S. M. Girvin and Terrence Jach, *Phys. Rev. B* **29**, 5617 (1984).
- [27] W. P. Su, Univ. of Illinois (Urbana) Preprint, 1985.
- [28] M. H. Friedman, J. B. Sokoloff, A. Widom, and Y. N. Srivastava, *Phys. Rev. Lett.* **52**, 1587 (1984).
- [29] R. B. Laughlin, private communications.
- [30] R. Tao, USC Preprint No. 085/003 (unpublished), 1985; R. Tao, *J. Phys. C* **18**, L1003 (1985).

Off-Diagonal Long-Range Order, Oblique Confinement, and the Fractional Quantum Hall Effect

S. M. Girvin

Surface Science Division, National Bureau of Standards, Gaithersburg, Maryland 20899

and

A. H. MacDonald

National Research Council, Ottawa, Ontario, Canada K1A OR6

(Received 24 November 1986)

We demonstrate the existence of a novel type of off-diagonal long-range order in the fractional-quantum-Hall-effect ground state. This is revealed for the case of fractional filling factor $v = 1/m$ by application of Wilczek's "anyon" gauge transformation to attach m quantized flux tubes to each particle. The binding of the zeros of the wave function to the particles in the fractional quantum Hall effect is a (2+1)-dimensional analog of *oblique confinement* in which a condensation occurs, not of ordinary particles, but rather of composite objects consisting of particles and gauge flux tubes.

PACS numbers: 72.20.My, 71.45.Gm, 73.40.Lq

A remarkable amount of progress has recently been made in our understanding of the fractional quantum Hall effect (FQHE)¹ following upon the seminal paper by Laughlin.² There remains, however, a major unsolved problem which centers on whether or not there exists an order parameter associated with some type of symmetry breaking.³⁻⁶ The apparent symmetry breaking associated with the discrete degeneracy of the ground state in the Landau gauge⁵ is an artifact of the toroidal geometry^{6,7} and is not an issue here. Rather, the questions that we are addressing have been motivated by the analogies which have been observed to exist^{4,8} between the FQHE and superfluidity and by recent progress towards a phenomenological Ginsburg-Landau picture of the FQHE.⁴ Further motivation has come from the development of the correlated-ring-exchange theory of Kivelson *et al.*⁹ (see also Chui, Hakim, and Ma,¹⁰ and Chui,¹⁰ and Baskaran¹¹). The existence of infinitely large ring exchanges is a signal of broken gauge symmetry in superfluid helium¹² and is suggestive of something similar in the FQHE. However, the concept of ring exchanges on large length scales has not as yet been fully

reconciled with Laughlin's (essentially exact⁷) variational wave functions which focus on the short-distance behavior of the two-particle correlation function. Furthermore it is clear that we cannot have an ordinary broken gauge symmetry since the particle density (which is conjugate to the phase) is ever more sharply defined as the length scale increases. The purpose of this Letter is to unify all these points of view by demonstrating the existence of a novel type of off-diagonal long-range order (ODLRO) in the FQHE ground state.

In second quantization the one-body density matrix is given by

$$\rho(z, z') = \sum_{m,n} \varphi_m^*(z) \varphi_n(z') \langle 0 | c_n^\dagger c_m | 0 \rangle, \quad (1)$$

where $\varphi_n(z)$ is the n th lowest-Landau-level single-particle orbital¹ in the symmetric gauge, and z is a complex representation of the particle position vector in units of the magnetic length.¹ It is an unusual feature of this problem that there is a unique single-particle state for each angular momentum. Hence by making only the assumption that the ground state is isotropic and homogeneous we may deduce $\langle 0 | c_n^\dagger c_m | 0 \rangle = v \delta_{nm}$, and thereby obtain the powerful identity:

$$\rho(z, z') = vg(z, z') = (v/2\pi) \exp(-\frac{1}{4} |z - z'|^2) \exp[\frac{1}{4} (z^* z' - z z^*)], \quad (2)$$

where $g(z, z')$ is the ordinary single-particle Green's function.¹³

We see from (2) that the density matrix is short ranged with a characteristic scale given by the magnetic length, just as occurs in superconducting films in a magnetic field.¹⁴ The same result can be obtained within first quantization via the expression

$$\rho(z, z') = \frac{N}{Z} \int d^2 z_2 \cdots d^2 z_N \Psi^*(z, z_2, \dots, z_N) \Psi(z', z_2, \dots, z_N), \quad (3)$$

where Z is the norm of Ψ .

If the lowest Landau level has filling factor $v = 1/m$ and the interaction is a short-ranged repulsion, then in the low-electron mass limit,⁷ the *exact* ground-state wave function is given by Laughlin's expression:

$$\Psi(z_1, \dots, z_N) = \prod_{i < j} (z_i - z_j)^m \exp\left(-\frac{1}{4} \sum_k |z_k|^2\right). \quad (4)$$

Laughlin's plasma analogy^{2,15} proves that the ground state is a liquid of uniform density so that Eq. (2) is valid. The rapid phase oscillations of the integrand in (3) cause ρ to be short ranged. There is, nevertheless, a peculiar type of long-range order hidden in the ground state. For reasons which will become clear below, this order is revealed by considering the singular gauge field used in the study of "anyons"^{16,17}:

$$\mathcal{A}_i(z_j) = \frac{\lambda\Phi_0}{2\pi} \sum_{i \neq j} \nabla_i \text{Im} \ln(z_j - z_i), \quad (5)$$

where $\Phi_0 = hc/e$ is the quantum of flux and λ is a constant. The addition of this vector potential to the Hamiltonian is not a true gauge transformation since a flux tube is attached to each particle. If, however, $\lambda = m$ where m is an integer, the net effect is just a change in the phase of the wave function

$$\Psi_{\text{new}} = \exp \left[-im \sum_{i < j} \text{Im} \ln(z_i - z_j) \right] \Psi_{\text{old}} \quad (6)$$

Application of (6) to the Laughlin wave function (4) yields

$$\begin{aligned} \Psi(z_1, \dots, z_N) \\ = \prod_{i < j} |z_i - z_j|^m \exp \left[-\frac{1}{4} \sum_k |z_k|^2 \right], \end{aligned} \quad (7)$$

which is purely real and is symmetric under particle exchange for both even and odd m . Hence we have the re-

$$\tilde{\rho}(z, z') = \frac{N}{Z} \int d^2 z_2 \cdots d^2 z_N \exp \left(-i \frac{e}{\hbar c} \int_z^{z'} d\mathbf{r} \cdot \mathcal{A}_1 \right) \Psi^*(z, z_2, \dots, z_N) \Psi(z', z_2, \dots, z_N), \quad (9)$$

where z and z' are vector representations of z and z' . The line integral in (9) is multiple valued but its exponential is single valued because the flux tubes are quantized. The additional phases introduced by the singular gauge transformation will cancel the phases in Ψ nearly everywhere, and produce ODLRO in $\tilde{\rho}$ if and only if the zeros of Ψ (which must necessarily be present because of the magnetic field¹⁹) are bound to the particles. Thus ODLRO in $\tilde{\rho}$ always signals a condensation of the zeros onto the particles (independent of whether or not the composite-particle occupation of the lowest momentum state diverges¹⁸). Because the gauge field \mathcal{A}_1 depends on the positions of all the particles, $\tilde{\rho}$ differs not just in the phase but in magnitude from ρ . This multiparticle object, which explicitly exhibits ODLRO, is very reminiscent of the topological order parameter in the XY model²⁰ and related gauge models^{21,22} and is intimately connected with the frustrated XY model which arises in the correlated-ring-exchange theory.⁹

For short-range interactions, the zeros of Ψ are directly on the particles and the associated phase factors are exactly canceled by the gauge term in (9) [see Eq. (7)]. As the range of the interaction increases, $m - 1$ of the zeros move away from the particles but remain nearby

markable result that both fermion and boson systems map into bosons in this singular gauge.

Substituting (7) into (3) and using Laughlin's plasma analogy,^{2,15} a little algebra shows that the singular-gauge density matrix $\tilde{\rho}$ can be expressed as

$$\tilde{\rho}(z, z')$$

$$= (v/2\pi) \exp[-\beta \Delta f(z, z')] |z - z'|^{-m/2}, \quad (8)$$

where $\beta \equiv 2/m$, and $\Delta f(z, z')$ is the difference in free energy between two impurities of charge $m/2$ (located at z and z') and a single impurity of charge m (with arbitrary location). Because of complete screening of the impurities by the plasma, the free-energy difference $\Delta f(z, z')$ rapidly approaches a constant as $|z - z'| \rightarrow \infty$. This proves the existence of ODLRO¹⁸ characterized by an exponent $\beta^{-1} = m/2$ equal to the plasma "temperature." For $m = 1$ the asymptotic value of Δf can be found exactly: $\beta \Delta f_\infty = -0.03942$. For other values of m , $\beta \Delta f(z, z')$ can be estimated either by use of the ion-disk approximation^{2,15} or the static (linear response) susceptibility of the (classical) plasma calculated from the known static structure factor⁸ (see Fig. 1).

The rigorous and quantitative results we have obtained above are valid for the case of short-range repulsive interactions for which Laughlin's wave function is exact. We now wish to use these results for a qualitative examination of more general cases and to deepen our understanding of the ODLRO. We begin by noting that $\tilde{\rho}$ can be rewritten in the ordinary gauge as

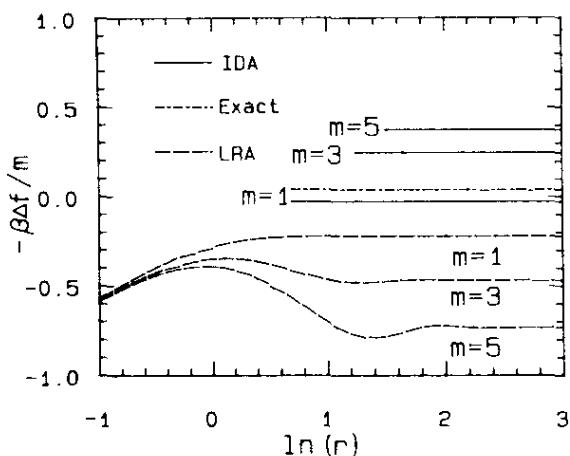


FIG. 1. Plot of $-\beta \Delta f(z, z')/m$ vs $r \equiv |z - z'|$ for filling factor $v = 1/m$. LRA is linear-response approximation, IDA is ion-disk approximation (shown only for separations exceeding the sum of the ion-disk radii). Because the plasma is strongly coupled, the IDA is quite accurate at $m = 1$ and improves further with increasing m . The LRA is less accurate at $m = 1$ and worsens with increasing m .

and bound to them.^{7,23} The gauge and wave function phase factors in (9) now appear in the form of the bound vortex-antivortex pairs. We expect such bound pairs *not* to destroy the ODLRO and speculate (based on our understanding of the Kosterlitz-Thouless transition²⁰) that the effect is at most to renormalize the exponent of the power law in (8). As the range of the potential is increased still further, numerical computations⁷ indicate that a critical point is reached at which the gap rather suddenly collapses and the overlap between Laughlin's state and the true ground state drops quickly to zero. We propose that this gap collapse corresponds to the unbinding of the vortices and hence to the loss of ODLRO and the onset of short-range behavior of $\tilde{\rho}(z, z')$. Recall that the distinguishing feature of the FQHE state is its long wavelength excitation gap. At least within the single-mode approximation,⁸ this gap can only exist when the ground state is homogeneous and the two-point correlation function exhibits perfect screening:

$$M_1 \equiv (v/2\pi) \int d^2r (r^2/2) [g(r) - 1] = -1.$$

In the analog plasma problem, the zeros of Ψ act like point charges seen by each particle and the M_1 sum rule implies that electrons see each other as charge- $(m = 1/v)$ objects; i.e., that m zeros are bound to each electron. Thus (within the single-mode approximation) there is a one-to-one correspondence between the existence of ODLRO and the occurrence of the FQHE.

The exact nature of the gap-collapse transition, which occurs when the range of the potential is increased,⁷ is not understood at present. However, it has been proven⁸ that the M_1 sum rule is satisfied by every homogeneous and isotropic state in the lowest Landau level. Hence the vortex unbinding should be a first-order transition to a state which breaks rotation symmetry (like the Tao-Thouless state²⁴) and/or translation symmetry (like the Wigner crystal^{4,8}). We know that as a function of temperature (for fixed interaction potential) there can be no

Kosterlitz-Thouless transition²⁰ since isolated vortices (quasiparticles) cost only a finite energy in this system^{4,25} (see, however, Ref. 10).

Further insight into the gap collapse can be obtained by considering the exceptional case of Laughlin's wave function with $m > 70$. In this case the zeros are still rigorously bound to the particles so that the analog plasma contains long-range forces (and $\tilde{\rho}$ exhibits ODLRO), but the plasma "temperature" has dropped below the freezing point.^{2,15} If such a state exhibits (sufficiently¹⁰) long-range positional correlations, the FQHE would be destroyed by a gapless Goldstone mode associated with the broken translation symmetry. Hence in this exceptional case the normal connection between ODLRO and the FQHE would be broken by a gap collapse due to positional order at a finite wave vector.

The existence of ODLRO in $\tilde{\rho}$ is the type of infrared property which suggests that a field-theoretic approach to the FQHE would be viable. It is clear from the results presented here that the binding of the zeros of Ψ to the particles can be viewed as a condensation,¹⁸ not of ordinary particles, but rather of composite objects consisting of a particle and m flux tubes. (We emphasize that these are *not* real flux tubes, but merely consequences of the singular gauge. The assumption that electrons can bind real flux tubes²⁶ is easily shown to be unphysical.²⁷) The analog of this result for hierarchical daughter states of the Laughlin states^{7,15} would be a condensation of composite objects consisting of n particles and m flux tubes (cf. Halperin's "pair" wave functions¹⁹). This seems closely analogous to the phenomenon of *oblique confinement*²² and it ought to be possible to derive the appropriate field theory from first principles by use of this idea.

Since the singular gauge maps the problem onto interacting bosons, coherent-state path integration²⁸ may prove useful. A step in this direction has been taken recently in the form of a Landau-Ginsburg theory which was developed on phenomenological grounds.⁴ In the static limit, the action has the "θ vacuum" form

$$S = \int d^2r [(-i\nabla + \mathbf{a})\psi(\mathbf{r})]^2 + i\phi(\mathbf{r})(\psi^*\psi - 1) - i(\theta/8\pi^2)(\phi\nabla \times \mathbf{a} + \mathbf{a} \times \nabla \phi), \quad (10)$$

where \mathbf{a} is not the physical vector potential but an effective gauge field⁴ representing frustration due to density deviations away from the quantized Laughlin density and ϕ is a scalar potential which couples both to the charge density and the "flux" density. From (10) the equation of motion for \mathbf{a} is (in the static case):

$$\theta \nabla \times \mathbf{a} = (\psi^*\psi - 1). \quad (11)$$

This equation and the parameter θ , which determines the charge carried by an isolated vortex, originally had to be chosen phenomenologically.⁴ Now, however, it can be justified by examination of Eq. (5) which shows that the curl of \mathcal{A}_I is proportional to the density of particles. If

we identify \mathbf{a} in (10) and (11) as

$$\mathbf{a} = \mathcal{A}_I + \mathbf{A}, \quad (12)$$

where \mathbf{A} is the physical vector potential and we take $\psi^*\psi$ as the particle density relative to the density in the Laughlin state, then Eq. (11) follows from (5) with the θ angle being given by $\theta = 2\pi/m$. This yields⁴ the correct charge of an isolated vortex (Laughlin quasiparticle) of $q^* = 1/m$. The connection between this result and the Berry phase argument of Arovas *et al.*²⁹ should be noted (see also Semenoff and Sodano³⁰). To summarize, it is the strong phase fluctuations induced by the frustration

associated with density variations [Eq. (11)] which pin the density at rational fractional values and account for the differences between the FQHE and ordinary superfluidity.⁴

We believe that these results shed considerable light on the FQHE, unify the different pictures of the effect, and emphasize the topological nature of the order in the zero-temperature state of the FQHE. The present picture leads to several predictions which can be tested by numerical computations by use of methods very similar to those now in use.³¹ ODLRO will be found only in states exhibiting an excitation gap. The decay of the singular-gauge density matrix will be controlled by the distribution of distances of the zeros of the wave function from the particles. This distribution, which can be artificially varied by changing the model interaction, directly determines the short-range behavior of the density-density correlation function and hence the ground-state energy.^{7,23}

The authors would like to express their thanks to C. Kallin, S. Kivelson, and R. Morf for useful conversations and suggestions.

¹The Quantum Hall Effect, edited by R. E. Prange and S. M. Girvin (Springer-Verlag, New York, 1986).

²R. B. Laughlin, Phys. Rev. Lett. **50**, 1395 (1983).

³P. W. Anderson, Phys. Rev. B **28**, 2264 (1983).

⁴S. M. Girvin, in Chap. 10 of Ref. 1.

⁵R. Tao and Yong-Shi Wu, Phys. Rev. B **30**, 1097 (1984).

⁶D. J. Thouless, Phys. Rev. B **31**, 8305 (1985).

⁷F. D. M. Haldane, in Chap. 8 of Ref. 1.

⁸S. M. Girvin, A. H. MacDonald, and P. M. Platzman, Phys. Rev. Lett. **54**, 581 (1985), and Phys. Rev. B **33**, 2481 (1986); S. M. Girvin in Chap. 9 of Ref. 1.

⁹S. Kivelson, C. Kallin, D. P. Arovas, and J. R. Schrieffer, Phys. Rev. Lett. **56**, 873 (1986).

¹⁰S. T. Chui, T. M. Hakim, and K. B. Ma, Phys. Rev. B **33**,

7110 (1986); S. T. Chui, unpublished.

¹¹G. Baskaran, Phys. Rev. Lett. **56**, 2716 (1986), and unpublished.

¹²R. P. Feynman, Phys. Rev. **91**, 1291 (1953).

¹³S. M. Girvin and T. Jach, Phys. Rev. B **29**, 5617 (1984).

¹⁴E. Brézin, D. R. Nelson, and A. Thiaville, Phys. Rev. B **31**, 7124 (1985).

¹⁵R. B. Laughlin, in Chap. 7 of Ref. 1.

¹⁶F. Wilczek, Phys. Rev. Lett. **49**, 957 (1982).

¹⁷D. P. Arovas, J. R. Schrieffer, F. Wilczek, and A. Zee, Nucl. Phys. B **251**, 117 (1985).

¹⁸We refer to this as ODLRO or condensation because of the slow power-law decay even though the largest eigenvalue $\lambda \equiv \int d^2z \hat{\rho}(z, z')$ of the density matrix diverges only for $m \leq 4$ [see C. N. Yang, Rev. Mod. Phys. **34**, 694 (1962)].

¹⁹B. I. Halperin, Helv. Phys. Acta **56**, 75 (1983).

²⁰J. V. José, L. P. Kadanoff, S. Kirkpatrick, and D. R. Nelson, Phys. Rev. B **16**, 1217 (1977).

²¹J. B. Kogut, Rev. Mod. Phys. **51**, 659 (1979).

²²J. L. Cardy and E. Rabinovici, Nucl. Phys. B **205**, 1 (1982); J. L. Cardy, Nucl. Phys. B **205**, 17 (1982).

²³D. J. Yoshioka, Phys. Rev. B **29**, 6833 (1984).

²⁴R. Tao and D. J. Thouless, Phys. Rev. B **28**, 1142 (1983).

The symmetric gauge version of this state exhibits threefold rotational symmetry.

²⁵A. M. Chang, in Chap. 6 in Ref. 1.

²⁶M. H. Friedman, J. B. Sokoloff, A. Widom, and Y. N. Srivastava, Phys. Rev. Lett. **52**, 1587 (1984), and **53**, 2592 (1984).

²⁷F. D. M. Haldane and L. Chen, Phys. Rev. Lett. **53**, 2591 (1984).

²⁸L. S. Schulman, Techniques and Applications of Path Integration (Wiley, New York, 1981).

²⁹D. Arovas, J. R. Schrieffer, and F. Wilczek, Phys. Rev. Lett. **53**, 722 (1984).

³⁰G. Semenoff and P. Sodano, Phys. Rev. Lett. **57**, 1195 (1986).

³¹F. D. M. Haldane and E. H. Rezayi, Phys. Rev. Lett. **54**, 237 (1985), and Phys. Rev. B **31**, 2529 (1985); F. C. Zhang, V. Z. Vulovic, Y. Guo, and S. Das Sarma, Phys. Rev. B **32**, 6920 (1985); G. Fano, F. Ortolani, and E. Colombo, Phys. Rev. B **34**, 2670 (1986).

Note added 1 H. Rezayi and F. D. M. Haldane have recently succeeded in computing the singular-gauge density matrix for a small number of particles on a sphere. The results are in complete accord with the discussion in the penultimate paragraph of this paper. See E. H. Rezayi and F. D. M. Haldane, Phys. Rev. Lett. Vol. **61** (1988) 1985.

Quantum Conductance in Networks

J. E. Avron, A. Raveh, and B. Zur

Department of Physics, Technion, Haifa 32000, Israel
(Received 21 November 1986)

We consider the quantum transport in networks. Arguments similar to those for the quantum Hall effect show that the averaged transport coefficients are quantized. Numerical calculations for a three-loop network yield the values 0, 1, and -1, depending on the fluxes threading the loops and the quantum state of the net. We characterize the conductance properties of such networks. We also discuss general properties of the transport coefficients in general multiloop networks.

PACS numbers 73.60.Aq, 02.40.+m, 72.20.My

It is known that there is a range of circumstances where the Hall conductance, at low temperatures, is a nonzero integer.¹ It is natural to inquire whether there are other systems with integer nonzero conductances. As we shall explain, networks are such systems: A network with L loops has $L(L-1)/2$ integer conductances which characterize the quantum state of the system and reflect its multiconnectivity. Like the Hall conductance, they are nondissipative and can be either holelike or electronlike, but unlike the Hall effect this does not reflect any band-structure properties.

The transport coefficients of the network, $2\pi g_{lm}$, are defined as the charge transported around loop l when the flux threading the m th loop, ϕ_m , increases adiabatically by 2π , the unit of quantum flux. Within linear-response theory, it turns out² that this is equivalent to the (time-averaged) ratio of the current in loop l to an infinitesimal emf acting on loop m .^{3,4} We shall concentrate on the cases where the network has three loops and l and m are distinct. We shall also assume throughout that the fluxes are changed sufficiently slowly for the adiabatic limit to hold. In particular the energy levels of the network are assumed to have nonvanishing gaps and we exclude situations where levels cross. Under these conditions, which guarantee no dissipation (dissipation arises when $l = m$ and the adiabatic limit does not hold), the nondiagonal conductances have nonlocal features. Also, the quantum (coherence) effects discussed below require temperatures which are low compared with a typical gap energy. Since energy gaps scale like (length)⁻² this favors small networks. This dictates temperatures in the millikelvin range and emf in microvolts for mesoscopic networks. Quantum coherence effects associated with the dissipative conductance in single mesoscopic loops, including nonlocal effects, are discussed by Sharvin and Sharvin.⁵ As yet, there are no experiments nor theory on the transport coefficients in two- or three-loop networks.

Consider, for example, a three-loop network made of mesoscopic, thin (metallic) wires (Fig. 1). Each loop is threaded by an independent flux tube ϕ_j , $j=1, 2, 3$. In comparison with the Hall effect, ϕ_3 plays the role of the magnetic field on the sample, ϕ_1 can be thought of as a

time-dependent flux replacing a battery, and ϕ_2 is the analog of Laughlin's flux tube. g_{12} is then the analog of the Hall conductance. Because of the analogy one may expect that g_{12} will be quantized and will be a nontrivial (antisymmetric) function of ϕ_1 . This, as we shall see, is essentially correct provided that suitable averaging is introduced: Let

$$\langle g_{lm} \rangle(\phi) \equiv \frac{1}{2\pi} \int_0^{2\pi} d\phi_m g_{lm}(\phi)$$

be the conductance averaged over the flux in the current loop m (ϕ denotes collectively the three fluxes). $\langle g_{12} \rangle(\phi_3)$ (or any other permutation of 1, 2, and 3) is an antisymmetric step-like function of ϕ_1 with steps at integer heights. This holds in great generality (i.e., even with electron-electron interaction, and also for thick wires and for more complicated networks) provided the system is in a pure quantum state which does not become degenerate as ϕ_1 and ϕ_2 are varied. It is a consequence of the fact that Kubo's formula for the (averaged) conductance has a topological interpretation being a first Chern number.^{2,6-10}

Here we shall describe parts of our numerical results and sketch the general theoretical structure. Details shall be presented elsewhere.¹¹

From the theory of superconducting networks¹² it is known that the analysis of the Schrödinger equation for the network of Fig. 1 (with one-dimensional wires) reduces to the study of 5×5 matrices (5 is the number of

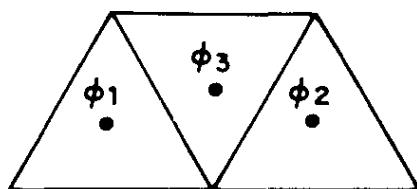


FIG. 1. Three-loop network with seven edges and five vertices. Each loop is threaded by a flux tube. The Hamiltonian for the network with point junctions is a 5×5 matrix. This network has nonzero quantized conductances.

vertices in the network) of the tight-binding type. It is therefore not surprising that the computation of the transport coefficients of the network also reduces to a (5×5) -matrix problem. The details of the reduction shall be given elsewhere.¹¹ In the matrix description the wave function sits on the vertices of the network. Consider the tight-binding Hamiltonian

$$H(v, v' | \phi) = n_v \delta_{vv} + (1 - \delta_{vv}) \sum_b [v, b] [v', b] \exp(-i[v, b] \gamma_b). \quad (1)$$

v and v' are vertex indices, n_v is the "coordination number" of the vertex v , b is a (directed) edge index, and $[v, b]$ is the incidence matrix, i.e., $[v, b] = 1$ if the edge b points into v , -1 if it points out of v , and 0 otherwise. $\gamma_b \equiv \int_b A$, where A is the vector potential associated to the fluxes ϕ . The lengths of all the edges b is set equal to one. $H(\phi)$ is identical to the de Gennes-Alexander¹² network Hamiltonian except on the diagonal. This slight modification makes it somewhat easier to handle. (The de Gennes-Alexander Hamiltonian gives an implicit eigenvalue problem.) Because of our interest in topological invariants the difference is presumably immaterial.

Diagonalizing the Hamiltonian one finds that the Chern number in the ground state, defined by Eq. (3) below, is

$$\langle g_{12} \rangle(\phi_1) = \begin{cases} 0 & \text{for } -\pi/3 < \phi_1 \bmod 2\pi < \pi/3, \\ 1 & \text{for } \pi/3 < \phi_1 \bmod 2\pi < \pi, \\ -1 & \text{for } -\pi < \phi_1 \bmod 2\pi < -\pi/3. \end{cases} \quad (2)$$

For the excited states one finds qualitatively similar, i.e., nontrivial, antisymmetric, periodic steplike functions that take the values 0, 1, and -1 . One also finds that $\langle g_{13} \rangle(\phi_2) = \langle g_{23} \rangle(\phi_1) = 0$ identically for all the states. Because of the topological nature of the results the fact that the network is made of three equilateral triangles is immaterial and one finds the same qualitative features in any network which is a deformation of Fig. 1.

To get a complete description and insight into the results we have to introduce some formalism. This is also necessary in order to describe the actual computation.

$H(\phi)$, the exact Schrödinger operator of the network, depends parametrically on the fluxes ϕ . For fixed ϕ , it has discrete spectrum. Because of the periodicity in the fluxes the parameter space can be identified with T^3 , the three-torus, i.e., we can identify ϕ_j with $\phi_j + 2\pi$.¹³

Let $P(\phi)$ denote a projection on a spectral subspace of $H(\phi)$ and C be a closed, two-dimensional surface in T^3 (equal to a closed two-chain). Suppose that $P(\phi)$ is smooth on C . It is a standard fact that the Chern number,^{2,6} $\text{Ch}(P, C) \equiv (i/2\pi) \int_C \text{Tr}[dP P dP]$, is an integer.

If the initial state of the system is given by $P(\phi)$ and there is no level crossing on $T_{lm}(\phi)$ (the two-dimensional slice of the three-torus going through ϕ and indexed by l and m), then Kubo's formula reads^{2,7-10} (see Ref. 10 for a rigorous derivation)

$$\langle g_{lm} \rangle(\phi) = \text{Ch}(P, T_{lm}(\phi)) \quad (3)$$

It follows¹¹ that $\langle g_{lm} \rangle(\phi)$ is gauge invariant, periodic in ϕ , antisymmetric in l and m , and independent of ϕ_l and

ϕ_m , and is quantized to be an integer. Also, in the absence of magnetic fields besides ϕ , which we shall assume, time reversal leads to the Onsager relation $\langle g_{lm} \rangle(\phi) = -\langle g_{lm} \rangle(-\phi)$.

It is known that the Chern numbers are closely related to degeneracies. Let D_q be the set of points where the q th gap in the energy spectrum closes. According to the von Neumann-Wigner theorem¹⁴ D_q is a discrete set. The second homology group of T^3/D_q is spanned by three two-tori, T_{12} , T_{23} , and T_{31} , and $|D_q|$ oriented two-spheres $S(\delta)$ that surround $\delta \in D_q$. An arbitrary closed two-chain in T^3/D_q can be written as a sum of the basic spheres and tori with integer coefficients. This relation lifts to a relation for the Chern numbers. It follows that the set of $3 + \sum |D_q|$ Chern numbers contains all the information about the net.

Relations among the basic Chern numbers follow from the following facts:

(1) $\sum_{\delta \in D_q} S(\delta)$ is homologous to zero, so that $\sum_{\delta \in D_q} \text{Ch}(P_q, S(\delta)) = 0$, for all q , where P_q is the projection on the spectral subspace with energies up to the q th level.

(2) The set D_q is invariant under inversion $\phi \rightarrow -\phi$, and for any δ in D_q , $\text{Ch}(P_q, S(\delta)) = \text{Ch}(P_q, S(-\delta))$.

For any closed two-chain c which is invariant under inversion, $\text{Ch}(P_q, c) = 0$. (4) If $H(\phi)$ is a periodic $n \times n$ matrix, P_n is the identity and so all its Chern numbers vanish.

For the Hamiltonian of Eq. (1) the set of points of degeneracy and their Chern numbers are given in Table I. Because of (2) above only points in the half-cube with $0 \leq \phi_j \leq \pi$ are listed. For the three basic tori we find

$$\begin{aligned} \text{Ch}(P_q, T_{23}(\phi_1)) &= \text{Ch}(P_q, T_{13}(\phi_2)) = 0, \\ \text{Ch}(P_q, T_{12}(\pi/2)) &= \delta(q, 1) \end{aligned} \quad (4)$$

This gives a complete characterization of the nondissipa-

TABLE I. Chern numbers for spheres surrounding points of degeneracy in the network of Fig. 1. $a \equiv \arccos(\frac{1}{2} - \sqrt{2})$.

Gap	Coordinates	Chern number
1	$(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})\pi$	1
1	$(\frac{4}{3}, \frac{1}{3}, 1)\pi$	-1
2	$(-a, a, 0)$	1
2	$(a, a, a/2)$	-1
3	$(\frac{1}{3}, \frac{4}{3}, \frac{1}{3})\pi$	1
3	$(\frac{1}{3}, \frac{4}{3}, 1)\pi$	-1

tive averaged conductance of the three-loop network of Eq. (1)

Equation (3), which relates the Chern numbers with the transport coefficients, is known to hold for the full Schrödinger equation of the network. We shall now describe how to extend this to matrix Hamiltonians. We shall consider here the case of matrix Hamiltonians

where the wave function sits on the bonds of the network. The case where it sits on vertices is more complicated and shall be dealt with elsewhere.¹¹

In a one-dimensional Schrödinger equation of a network the wave function is $\Psi_b(x_b)$, $0 \leq x_b \leq 1$. The coordinate x_b is measured in the b direction. On b , Ψ solves the one-dimensional Schrödinger equation in a gauge potential and therefore has the form

$$\Psi_b(x) = [T(\psi_+, \psi_-)](x) \equiv \exp[i\gamma_b(x)][\psi_+(b)\exp(ikx) + \psi_-(b)\exp(-ikx)], \quad (5)$$

where $\gamma_b(x) \equiv \int^x A$. $\psi_+(b)$ and $\psi_-(b)$ denote the amplitudes of the forward- and backward-moving waves on the bond b . The linear operator T is introduced for later purposes. A vertex with n_v edges connected to it is described by a unitary $n_v \times n_v$ scattering matrix, S_v ,¹⁵ which maps the incoming waves on the outgoing waves, $S_v \Psi_{v,\text{in}} = \Psi_{v,\text{out}}$. $\Psi_{v,\text{in/out}}$ are n_v -dimensional vectors of complex numbers:

$$\begin{aligned} \Psi_{v,\text{in}}(b) &\equiv \delta(1, [v, b]) \psi_+(b) \exp[ik + i\gamma_b] + \delta(-1, [v, b]) \psi_-(b), \\ \Psi_{v,\text{out}}(b) &\equiv \delta(1, [v, b]) \psi_-(b) \exp[-ik + i\gamma_b] + \delta(-1, [v, b]) \psi_+(b). \end{aligned} \quad (6)$$

b runs over the n_v edges associated to v . The unitarity of S_v guarantees that current is conserved at each vertex.

A basic tool is this: Consider $T: H_1 \rightarrow H_2$, where $H_{1,2}$ are two Hilbert spaces (not necessarily of the same dimension). Let Q be an orthogonal projection on H_2 , and suppose that $TT^*Q = Q$ (T^* denotes the adjoint of T). Suppose that $QdT T^*$ is smooth and globally defined. Then $P \equiv T^*QT$ is a projection on H_1 and the Chern numbers for the two projections coincide.

We apply this to H_1 , the finite-dimensional complex vector space, and H_2 , the Hilbert space of functions. The map T from C^2 to $L^2[0,1]$ is given by Eq. (5). The scalar product in C^2 , induced by the scalar product in L^2 , has the "Riemann metric" A where

$$A_{ij} \equiv \delta_{ij} + [(1 - \delta_{ij})/k] \exp(-ik\epsilon_{ij}) \sin(k), \quad i, j = 1, 2$$

ϵ_{ij} is the completely antisymmetric tensor. The metric is independent of the gauge field and is nonsingular provided $k \neq 0$. One finds¹¹ for dT^*T

$$dT^*T = d(A^{-1})^*A - i(A^{-1})^* \left[\begin{array}{cc} \int [d\gamma(x) + x dk] dx & \int [d\gamma(x) - x dk] \exp(-2ikx) dx \\ \int [d\gamma(x) + x dk] \exp(2ikx) dx & \int [d\gamma(x) - x dk] dx \end{array} \right]. \quad (7)$$

γ is linear in ϕ , so that $d\gamma$ is independent of ϕ . Q , k , A^{-1} , and $QdT T^*$ are all smooth in ϕ provided no levels cross and k does not vanish. This establishes the equality of the Chern numbers.

In summary, networks have quantized averaged conductances which are nontrivial in networks with three or more loops. The computation of Chern numbers for network Hamiltonians with one-dimensional connecting links reduces to the study of finite matrices. Finally, homology provides a convenient and compact way of presenting the conductance functions of networks.

This work was supported by the Israel Academy of Sciences, Minerva, and the U.S.-Israel Binational Science Foundation Grant No. 84-00376. We are indebted to A. Libchaber for a conversation that started this work, to Y. Imry, S. Lipson, R. Seiler, S. Shtrikman, and J. Zak for discussions, and to D. Vollhardt for help with the references.

45, 494 (1980), H. L. Stormer *et al.* Phys. Rev. Lett. **56**, 85 (1986).

²J. E. Avron and R. Seiler, Phys. Rev. Lett. **54**, 259 (1985).

³Loop currents I_e , $e \in E$, are related to the usual edge currents I_e , $e \in E$, by $I_e = \sum_{f \in F} [e, f] I_{ef}$. $[e, f]$ is the incidence matrix of the graph.

⁴R. B. Laughlin, Phys. Rev. B **23**, 5632 (1981).

⁵D. Yu. Sharvin and Yu. V. Sharvin, Pis'ma Zh. Teksp. Teor. Fiz. **34**, 285 (1981) [JETP Lett. **34**, 272 (1981)]; R. A. Webb *et al.*, Phys. Rev. Lett. **51**, 690 (1983); M. Büttiker *et al.*, Phys. Lett. **96A**, 365 (1983); C. P. Umbach *et al.*, Phys. Rev. B **30**, 4048 (1984); Y. Gefen, Y. Imry, and M. Azbel, Phys. Rev. Lett. **52**, 129 (1984); V. Chandrasekhar *et al.*, Phys. Rev. Lett. **55**, 1610 (1985); R. A. Webb *et al.*, Phys. Rev. Lett. **54**, 2696 (1985); M. Büttiker *et al.*, Phys. Rev. B **31**, 6207 (1985); M. Büttiker, Phys. Rev. Lett. **57**, 1761 (1986); A. D. Benoit *et al.*, Phys. Rev. Lett. **56**, 1765 (1986).

⁶D. J. Thouless, M. Kohmoto, M. P. Nightingale, and M. den Nijs, Phys. Rev. Lett. **49**, 405 (1982).

⁷Q. Niu and D. J. Thouless, J. Phys. A **17**, 2453 (1984).

⁸Q. Niu, D. J. Thouless, and Y. S. Wu, Phys. Rev. B **31**,

¹K. von Klitzing, G. Dorda, and M. Pepper, Phys. Rev. Lett.

3372 (1985).

⁹R. Tao and F. D. M. Haldane, Phys. Rev. B **33**, 3844 (1986).

¹⁰J. E. Avron, R. Seiler, and L. Yaffe, Commun. Math. Phys. (to be published).

¹¹J. E. Avron, A. Raveh, and B. Zur, to be published.

¹²P. G. de Gennes, C. R. Acad. Sci., Ser. B **292**, 9, 279 (1981); R. Rammal and G. Toulouse, Phys. Rev. Lett. **49**,

1194 (1982); S. Alexander, Phys. Rev. B **27**, 1541 (1983); E. Domany *et al.*, Phys. Rev. B **26**, 3110 (1983).

¹³The canonical choice of $H(\phi)$ is not periodic in ϕ . There is, however, a choice that makes $H(\phi)$ periodic. See Ref. 7 for more details.

¹⁴J. von Neumann and E. Wigner, Phys. Z. **30**, 467 (1929).

¹⁵P. W. Anderson *et al.*, Phys. Rev. B **22**, 3519 (1980); B. Shapiro, Phys. Rev. Lett. **50**, 747 (1983).

Superconducting Ground State of Noninteracting Particles Obeying Fractional Statistics

R. B. Laughlin

*Department of Physics, Stanford University, Stanford, California 94305, and
University of California, Lawrence Livermore National Laboratory, Livermore, California 94550*

(Received 28 January 1988)

In a previous paper, Kalmeyer and Laughlin argued that the elementary excitations of the original Anderson resonating-valence-bond model might obey fractional statistics. In this paper, it is shown that an ideal gas of such particles is a new kind of superconductor.

PACS numbers: 74.65.+n, 05.30.-d, 67.40.-w, 75.10.Jm

In a recent Letter,¹ Kalmeyer and I proposed that the ground state of the frustrated Heisenberg antiferromagnet in two dimensions and the fractional quantum Hall state for bosons might be the "same," in the sense that the two systems could be adiabatically evolved into one another without crossing a phase boundary. Whether or not this is the case is not presently clear. Indeed, the existence of a spin-liquid state of *any* spin- $\frac{1}{2}$ antiferromagnet in two dimensions has not been demonstrated. However, the case for a phase boundary's not being crossed is sufficiently strong that it is appropriate to ask what the consequences would be if this occurred. Adiabatic evolution is a particularly useful concept in the study of fractional quantum Hall "matter." So long as the energy gap remains intact, the "charge" of its fractionally charged excitations remains *exact* and the concomitant long-range forces between them, their fractional statistics, remain operative. This is why the fractional quantum Hall effect is so stable and reproducible. The persistence of the gap under evolution of the fractional quantum Hall problem into the magnet problem would allow us to make exact statements about the magnet without knowing *anything* about its Hamiltonian. In particular, the excitation spectrum of the magnet would be almost identical to that proposed by Kivelson, Rokhsar, and Sethna,² and completely within the spirit of the Anderson resonating-valence-bond idea,^{3,4} except for one crucial detail: Both the chargeless spin- $\frac{1}{2}$ excitations, the "spinons," and the charged spinless excitations, the "holons," would obey $\frac{1}{2}$ fractional statistics.^{5,6} The purpose of this Letter is to point out that this overlooked property may well account for high-temperature superconductivity.

Kalmeyer and I found the magnetic analog of the charge- $\frac{1}{2}$ quasiparticle of the fractional quantum Hall effect to be a spin- $\frac{1}{2}$ excitation, well described qualitatively as a spin-down electron on site j surrounded by an otherwise featureless spin liquid. This particle is our version of the "spinon." Like the quasiparticle of the fractional quantum Hall state, it carries a "charge," that is, its spin, that is in a deep and fundamental sense fractional. In the limit that the antiferromagnetic interactions are turned off, the excitation spectrum of the magnet is

purely bosonic. Spin- $\frac{1}{2}$ particles occur because these "elementary" excitations are fractionalized: Half the boson is deposited in the sample interior and half at the boundary. It was first pointed out by Halperin⁶ that, in the fractional quantum Hall effect, the fractionalization of the electron charge e into the quasiparticle charge $\frac{1}{3}e$ causes the quasiparticle to obey $\frac{1}{3}$ fractional statics. That is, each quasiparticle acts as though it were a boson carrying a magnetic solenoid containing magnetic flux $\frac{1}{3} \times hc/e$. This fact, deduced by Halperin from the experimentally observed fractional quantum Hall hierarchical states, was later shown by me⁷ to follow from the analytic properties of the quasiparticle wave functions. It arises physically because the states available to the multiquasiparticle system must be enumerated differently from those available to fermions or bosons. In other words, it comes from counting. Now, it is clear by inspection that the preferred nature of this representation does not care about the existence of a lattice. Thus the validity of our identification clearly predicts that spinons obey $\frac{1}{2}$ statistics.

Let us now imagine doping this lattice with holes. The most natural way to do this, in my opinion, is first to make a spinon, thus fixing the spin on site j , and then remove the electron possessing that spin. It is necessary to make the spinon first because an electron cannot be removed before its spin state is known. If one simply rips an "up" electron from site j , one tacitly projects the ground state onto the set of states with the j th spin up, thus creating an excitation with spin 1. This may be thought of as a pair of spinons in close proximity. Unless the interaction between spinons is attractive and sufficiently large (Kalmeyer and I found it to be repulsive¹), to make this "spin wave" will be more expensive energetically than to make an isolated spinon. Given that this occurs, the resulting spinless particle, the "holon," should also exhibit $\frac{1}{2}$ fractional statistics because it is a composite of a spinon and a fermion.

Assume now that we have a gas of such holons obeying fractional statistics. What are its properties expected to be? This question was addressed to some extent by Arovas *et al.*,⁸ who computed the second virial coefficient of an ideal gas of particles obeying fractional statistics as

a function of the fraction v . Not surprisingly, they found a smooth interpolation between the case of fermions, which acts like a classical gas with *repulsive* interactions, and that of bosons, which acts like a classical gas with *attractive* interactions. Thus, if we insist on thinking of these particles as fermions, we must conclude that there is an enormous attractive force between them. This is also evident when one considers the low-temperature properties. Fermions at density ρ have a large degeneracy pressure, and thus a large internal energy, while bosons have neither. Since fractional-statistics particles are in between, they have, *vis-à-vis* fermions, attractive forces comparable in scale to the Fermi energy. It is also important that spinless particles obeying fractional statistics cannot undergo Bose condensation. They are not bosons. However, if the fraction is $\frac{1}{2}$, then *pairs* of particles are bosons.

There is therefore good reason to suspect that a gas of particles obeying $\frac{1}{2}$ statistics might actually be a superconductor with a charge-2 order parameter. Let us investigate this possibility by considering a gas of fractional-statistics particles described by the free-particle Hamiltonian

$$\mathcal{H} = \sum_j \frac{p_j^2}{2m} \quad (1)$$

Any eigenstate of this Hamiltonian may be written in the manner

$$\begin{aligned} \Psi(z_1, \dots, z_N) \\ = \left[\prod_{j < k} \frac{(z_j - z_k)^v}{|z_j - z_k|^1} \right] \Phi(z_1, \dots, z_N), \end{aligned} \quad (2)$$

where z_j denotes the position of the j th particle in the x - y plane expressed as a complex number, $v = \frac{1}{2}$, and Φ is a Fermi wave function. This is the singular gauge transformation first discussed by Wilczek.⁵ If we have an eigenstate Ψ satisfying $\mathcal{H}\Psi = E\Psi$, then Φ satisfies

$$\mathcal{H}_{HF} = \frac{1}{2} E_0 + (-E_0 - \frac{1}{4}) \Pi_0 + \sum_{n=0}^{\infty} \left[n + \sum_{k=1}^n (-1)^k \binom{n}{k} \left(\frac{1}{4} \sum_{l=1}^k \frac{1}{l} - \frac{1}{2k} \right) + \frac{1}{4(n+1)} \right] \Pi_n, \quad (6)$$

with

$$E_0 = \int_0^\infty r^{-1} [1 - e^{-r^2/\alpha^2}] e^{-ar} dr, \quad (7)$$

in units of the equivalent cyclotron frequency $\hbar\omega_c = 2\pi v(\hbar^2/m)\rho$, where Π_n denotes the projector onto the n th Landau level, and α is a regulation parameter, effectively the inverse of the sample radius. Since \mathcal{H}_{HF} preserves Landau-level index, the state we guessed is a true variational minimum. Note, however, the logarithmic divergence in the Lagrange-multiplier spectrum, im-

$\mathcal{H}'\Phi = E\Phi$ where

$$\mathcal{H}' = \sum_j \frac{1}{2m} |\mathbf{p}_j + \mathbf{A}_j|^2, \quad (3)$$

and

$$\mathbf{A}_j(\mathbf{r}_j) = v \sum_{k \neq j} \hat{\mathbf{z}} \times \mathbf{r}_{jk} / |\mathbf{r}_{jk}|^2. \quad (4)$$

Thus, in the Fermi representation, each particle appears to carry a magnetic solenoid with it as it moves around in the sample. The vector potential felt by a particle is then the sum of the vector potentials generated by all the other particles. Because particles obeying $\frac{1}{2}$ statistics behave like fermions, in the sense that they possess degeneracy pressure, let us attempt to solve this problem in the Hartree-Fock approximation: We make a variational wave function that is a single Slater determinant constructed of orbitals $\phi_j(z)$ and minimize the expected energy. The orbitals then obey equations of the form

$$\mathcal{H}_{HF}\phi_j(z) = \lambda_j \phi_j(z), \quad (5)$$

where \mathcal{H}_{HF} is the first variation of $\langle \mathcal{H}' \rangle$ and λ_j is a Lagrange multiplier. The latter has the physical sense of a partial derivative of the total energy with respect to occupancy of the j th orbital. Since, in the mean-field sense, each particle must see a uniform density of magnetic solenoids carrying flux vhc/e , it is reasonable to guess the solution to be Landau levels, with the magnetic length a_0 related to the particle density ρ by $a_0^2 = (2\pi v\rho)^{-1}$. Self-consistency is achieved when the lowest $1/v$ Landau levels are filled. Thus, the fractions $v = 1, \frac{1}{2}, \frac{1}{3}, \dots$ are special cases in which a gap opens up in the fermionic spectrum.

Let us now test these equations in a case for which we know the answer, namely $v = 1$, the noninteracting Bose gas. If the variational procedure describes this limit correctly, there is good reason to trust its predictions for $v = \frac{1}{2}$. Evaluating the self-consistent field with one Landau level filled, I obtain

plying that the cost to inject either a "particle" or an "antiparticle" is arbitrarily large. This is absolutely the correct result. The noninteracting Bose gas has no low-lying fermionic excitations. The fact that these divergences are logarithmic suggests that the relevant excitations are actually quantum vortices. That this is, in fact, the case may be seen by our imagining an extra particle to be placed at the origin and calculating the expected current density $\langle \mathbf{J}(r) \rangle$. The current-density operator may be written $\mathbf{J}(r) = m^{-1} (\mathbf{p} + \mathbf{A}_{old} + \Delta\mathbf{A})$, where \mathbf{A}_{old} is the vector potential in the absence of the extra particle

and $\Delta\mathbf{A}$ is the vector potential generated by a solenoid at the origin. Since $(\mathbf{p} + \mathbf{A}_{\text{old}}) = 0$, the current density must just be the particle density at r times $\Delta\mathbf{A}$, or a vortex of magnetic strength hc/e .

The expected energy of the ground state state is $N/4$ in these units. This is considerably higher than the correct answer of zero. This discrepancy is due to the fact that the wave function is forced by its construction to go to zero when the particles come together. It is thus more appropriate for the description of real helium than noninteracting bosons. It should also be noted that this behavior is actually required of the $v = \frac{1}{2}$ wave function. Let us observe finally that the broken symmetry characteristic of a superfluid is not expressly exhibited by the Hartree-Fock ground state. This is as expected. It was shown by Bogoliubov⁹ that the broken symmetry of a Bose gas is absent unless the bosons interact. All that is required for the symmetry to break is a weak interparti-

cle repulsion and the presence in the "unperturbed" Bose gas of a collective mode dispersing quadratically with the mass of the bare particles. In the present case, it is easy to see that the variational solution possesses a collective mode that disperses quadratically. Since $\langle \mathcal{H}' \rangle / N$ is proportional to the particle density, the pressure is constant, and thus the bulk modulus is zero. It is a straightforward matter to calculate the mass of this mode by the magnetoexciton procedure of Kallin and Halperin.¹⁰ My preliminary results give a value of approximately $\frac{1}{2}$ the bare mass. The precise value of this mass is not so important as the fact that it is of order unity. The collective mode may be thought of both as a density wave and as a magnetoexciton consisting of a hole in the lowest Landau level and a particle in the first excited Landau level, bound together by a logarithmic potential.

Let us now turn to the case of interest, $v = \frac{1}{2}$. It is so similar to the $v = 1$ case that there is little to say. Assuming two Landau levels filled, I obtain

$$\begin{aligned} \mathcal{H}_{HF} = & \frac{11}{16} + \frac{1}{4} E_0 + \left(-\frac{1}{2} E_0 - \frac{1}{8} \right) \Pi_0 + \left(-\frac{1}{2} E_0 + \frac{29}{24} \right) \Pi_1 \\ & + \sum_{n \geq 2} \left[n + \sum_{k=1}^n \binom{n}{k} (-1)^k \left(\frac{1-k}{4} \sum_{l=1}^k \frac{1}{l} - \frac{1}{4k} \right) + \frac{3}{8(n+1)} - \frac{1}{8(n+2)} \right] \Pi_n. \quad (8) \end{aligned}$$

Thus, we again have a true variational solution with vortexlike fermionic excitations. Repeating the arguments for $v = 1$, I find that the flux quantum to which the vortices correspond is $hc/2e$, exactly as expected of a charge-2 superfluid. Once again, a soft collective mode will mix into the ground state to break the symmetry when repulsive interactions are introduced. Thus, the ground state is a superfluid very similar to liquid helium except that the charge of its order parameter is 2.

While considerable work needs to be done to quantify this picture, some of its implications may be seen at a glance. By far the most important is that a normal-metal state, in the sense of Fermi-liquid theory, does not exist, just as Anderson⁴ suggested. A corollary is that the occurrence of superconductivity does not have anything to do with self-consistent opening of an energy gap in the tunneling spectrum, as occurs in the BCS theory. Indeed, I find that tunneling cannot even be understood outside the context of the creation of spinons by the tunneling event. It should be noted that this is also consistent with Anderson's views.¹¹ A critical prediction is that an energy gap *must* occur in the spin-wave spectrum, the spin analog of the collective mode¹² of the fractional quantum Hall state. This is because the presence or absence of this gap is precisely the difference between the disordered and ordered states.

In summary, it is possible that high- T_c superconductivity can be accounted for by the following simple idea: The force mediated by the spins of the Mott insulator is

not an attractive potential, but rather an attractive *vector* potential.

I gratefully acknowledge numerous helpful conversations with S. Kivelson, J. Sethna, V. Kalmyer, L. Susskind, A. L. Fetter, P. W. Anderson, F. Wilczek, B. I. Halperin, J. R. Schrieffer, T. H. Geballe, M. R. Beasley, and A. Kapitulnik. This work was supported primarily by the National Science Foundation under Grant No. DMR-85-10062 and by the National Science Foundation-Materials Research Laboratory Program through the Center for Materials Research at Stanford University. Additional support was provided by the U.S. Department of Energy through the Lawrence Livermore National Laboratory under Contract No. W-7405-Eng-48.

¹V. Kalmyer and R. B. Laughlin, Phys. Rev. Lett. **59**, 2095 (1987)

²S. Kivelson, D. Rokhsar, and J. Sethna, Phys. Rev. B **35**, 8865 (1987)

³P. W. Anderson, Mater. Res. Bull. **8**, 153 (1973)

⁴P. W. Anderson, Science **236**, 1196 (1987), P. W. Anderson, G. Baskaran, Z. Zou, and T. Hsu, Phys. Rev. Lett. **58**, 2790 (1987).

⁵F. Wilczek, Phys. Rev. Lett. **49**, 957 (1982), F. Wilczek and A. Zee, Phys. Rev. Lett. **51**, 2250 (1983)

⁶B. I. Halperin, Phys. Rev. Lett. **52**, 1583 (1984)

⁷R. B. Laughlin, in *The Quantum Hall Effect*, edited by

R. E. Prange and S. M. Girvin (Springer-Verlag, New York, 1987), p. 233.

⁸D. P. Arovas, R. Schrieffer, F. Wilczek, and A. Zee, Nucl. Phys. **B251** [FS13], 117 (1985).

⁹N. N. Bogoliubov, J. Phys. (Moscow) **11**, 23 (1947); a good discussion of this may be found in A. L. Fetter and J. D.

Walecka, *Quantum Theory of Many-Particle Systems* (McGraw-Hill, New York, 1971), p. 313.

¹⁰C. Kallin and B. I. Halperin, Phys. Rev. B **30**, 5655 (1984).

¹¹P. W. Anderson, private communication.

¹²S. M. Girvin, A. H. MacDonald, and P. M. Platzman, Phys. Rev. Lett. **54**, 581 (1985).

An example of phase holonomy in WKB theory

Michael Wilkinson

H H Wills Physics Laboratory, Royal Fort, Tyndall Avenue, Bristol, BS8 1TL, UK†

Received 4 June 1984

Abstract. This paper discusses the application of WKB theory to Harper's equation

$$\psi_{n+1} + \psi_{n-1} + 2\alpha \cos(2\pi\beta n + \delta) \psi_n = E\psi_n$$

in the case in which β is very close to a rational number, p/q .

The WKB wavefunction for this system is a vector valued quantity, proportional to an eigenvector \mathbf{u} of a matrix $\hat{H}(x, p)$, which is parametrised by the phase space coordinates x and p . The complex phase of \mathbf{u} is determined by a non-holonomic connection rule; when transported around a cycle and in phase space, \mathbf{u} is multiplied by a phase factor $e^{i\gamma_c}$. This phase change manifests itself as a modification of the Bohr-Sommerfeld quantisation condition.

1. Introduction

This paper describes an unusual form of Bohr-Sommerfeld quantisation, involving a holonomy argument. As well as being interesting in its own right, the method discussed here can be applied to the difficult problem of finding the Bohr-Sommerfeld quantisation condition for Bloch electrons in a magnetic field. The system treated in this paper is a simplified model for this problem, which is often called Harper's equation. This model will be introduced in § 2; the remainder of this introduction will describe the principle of the method.

For the system considered, the WKB wavefunction can be thought of as a vector-valued quantity, given by

$$\psi(x) = A(x)\mathbf{u}(x) \exp\left(\frac{i}{\hbar} \int^x p(x') dx'\right) \quad (1.1)$$

where $p(x)$ and $A(x)$ are slowly varying functions, and the vector \mathbf{u} is a solution of the eigenvalue equation

$$\hat{H}(x, p)\mathbf{u} = \epsilon\mathbf{u}. \quad (1.2)$$

In equation (1.2), ϵ is the energy of the solution $\psi(x)$, and \hat{H} a complex Hermitian matrix which is a function of two parameters x and p . Since the energy $E = \epsilon(x, p)$ is a constant for a given solution, equation (1.2) gives both \mathbf{u} and p as functions of x , as in (1.1). The curves in the $x-p$ plane defined by $E = \epsilon(x, p) = \text{constant}$ are called phase trajectories. When the phase trajectories given by (1.2) are closed orbits, then a solution $\psi(x)$ must remain single-valued when it is traced around the phase trajectory.

† Address after September 1st 1984: Department of Physics, California Institute of Technology, Pasadena, California 91125, USA

This condition is only satisfied for certain values of E , which are determined by a Bohr–Sommerfeld quantisation condition.

In equation (1.2), the eigenvector $\mathbf{u}(x, p)$ is determined only up to a complex-valued multiplying constant, or, if \mathbf{u} is assumed to be normalised, up to a complex phase factor $e^{i\theta}$. This phase factor can be determined by requiring that the amplitude $A(x)$ in (1.1) be real. Given this condition on $A(x)$, it will be shown how the WKB theory for the system leads to a connection formula, by means of which the vector \mathbf{u} can be transported through the phase space with its phase fully determined. It turns out that this phase connection is non-holonomic, so that when \mathbf{u} is transported clockwise around a closed circuit in phase space, it is multiplied by a phase factor $e^{i\gamma}$.

This phase factor affects the Bohr–Sommerfeld quantisation condition. Consider the phase change of the solution (1.1) after making one circuit of a closed phase trajectory. This has contributions from the oscillatory term, from a pair of turning points where $p(x) = 0$ and $A(x)$ diverges, plus a contribution $\gamma(E)$ from the phase factor evaluated for a phase trajectory of energy E . The condition for the wavefunction ψ to be single valued is therefore

$$2\pi n = \frac{1}{\hbar} \oint_{E=E_n} p(x) dx + \frac{\pi}{2} + \frac{\pi}{2} + \gamma(E_n), \quad (1.3)$$

or

$$\oint_{E=E_n} p(x) dx = [2\pi(n + \frac{1}{2}) - \gamma(E_n)]\hbar. \quad (1.4)$$

This equation (1.4) is the Bohr–Sommerfeld quantisation condition determining the eigenvalues E_n of the system.

The plan of this paper is as follows. Section 2 introduces the system under consideration and discusses how WKB theory can be applied to this system. Section 3 derives an asymptotic formula for the product of a string of slowly varying transfer matrices. Section 4 applies this formula to the WKB problem for Harper's equation, and § 5 obtains the Bohr–Sommerfeld quantisation condition. Section 6 summarises the theoretical results and compares them with numerical values, and § 7 discusses the connections between this work and recent work on adiabatic theory and the quantised Hall effect.

2. WKB analysis of Harper's equation

The system analysed in this paper is Harper's equation

$$\psi_{n+1} + \psi_{n-1} + 2\alpha \cos(2\pi\beta n + \delta)\psi_n = E\psi_n, \quad (2.1)$$

which is frequently used in models for Bloch electrons in a magnetic field, and as a model for electrons in an incommensurate potential (Harper 1955, Simon 1982). As pointed out by Sokoloff (1981), solutions of (2.1) can be obtained by a WKB method whenever β is sufficiently close to a rational number, p/q (where p and q are coprime integers). The condition for WKB theory to be applicable is

$$|q^2\Delta\beta| \ll 1, \quad \Delta\beta = \beta - p/q, \quad (2.2)$$

and for almost all β , there exist values of p/q for which $|q^2\Delta\beta|$ is arbitrarily small. This follows from a property of continued fractions (Khinchin 1964).

Before describing how WKB theory can be applied to (2.2), it will be useful to consider the case $\beta = p/q$, so that the coefficients of the difference equation (2.1) are periodic with period q . In this case, therefore, Bloch's theorem applies and exact solutions can be obtained; the Bloch solution has the form

$$\psi_n = e^{ikn} u_n(\delta, k), \quad (2.3)$$

where u_n is periodic with period q

$$u_{n+q} = u_n. \quad (2.4)$$

This result can also be written in terms of a set of Fourier amplitudes for u_n

$$\psi_n = e^{ikn} \sum_{m=0}^{q-1} a_m \exp(2\pi i pmn/q), \quad (2.5)$$

which will prove more useful for some purposes.

From equation (2.1), it can be seen that the q -component vectors u_n or a_m can be determined as eigenvectors of a $q \times q$ complex Hermitian matrix. To distinguish these matrices and vectors from some two-component vectors and 2×2 matrices which will be introduced later, a quantum mechanical notation will be used;

$$\hat{H}(\delta, k)|u(\delta, k)\rangle = \epsilon|u(\delta, k)\rangle, \quad (2.6)$$

where the matrix elements of \hat{H} are given by

$$\begin{pmatrix} 2\alpha \cos(2\pi\beta + \delta) & e^{ik} & & & & \\ \vdots & \ddots & \ddots & & & \\ & \ddots & e^{-ik} & 2\alpha \cos(2\pi\beta n + \delta) & & \\ & & 0 & \ddots & \ddots & \\ & & & \ddots & e^{ik} & 2\alpha \cos(2\pi\beta q + \delta) \\ e^{ik} & & & & \ddots & \end{pmatrix} \begin{pmatrix} e^{-ik} \\ u_1 \\ \vdots \\ u_n \\ \vdots \\ u_q \end{pmatrix} = \epsilon \begin{pmatrix} u_1 \\ \vdots \\ u_n \\ \vdots \\ u_q \end{pmatrix}. \quad (2.7)$$

It is fairly easy to show that the eigenvalues E are periodic in both δ and k with period $2\pi/q$; in fact the q eigenvalues are given by the equation

$$f(E) = \cos qk + \alpha^q \cos q\delta, \quad (2.8)$$

where f is a q th degree polynomial (Wilkinson 1984). The q different sheets of $\epsilon(\delta, k)$ normally do not touch each other. When q is even, however, one pair of sheets of $\epsilon(\delta, k)$ does touch at isolated points in the $\delta - k$ plane (Bellissard and Simon 1982).

There is another, complementary, method for analysing equation (2.1) when β is rational (i.e. $\Delta\beta = 0$); this is the transfer matrix method. It is easy to see that equation (2.1) can be written in the form

$$\begin{pmatrix} \psi_{n+1} \\ \psi_n \end{pmatrix} = \tilde{T}(x_n, E) \begin{pmatrix} \psi_n \\ \psi_{n-1} \end{pmatrix}, \quad (2.9)$$

where

$$\tilde{T}(x, E) = \begin{pmatrix} E - 2\alpha \cos x & -1 \\ 1 & 0 \end{pmatrix},$$

$$x_n = 2\pi\beta n + \delta. \quad (2.10)$$

3462

M Wilkinson

Consider a transfer matrix $\tilde{M}(x, E)$ describing a 'jump' of q steps

$$\begin{pmatrix} \psi_{(n+1)q+1} \\ \psi_{(n+1)q} \end{pmatrix} = \tilde{M}(x_n, E) \begin{pmatrix} \psi_{nq+1} \\ \psi_{nq} \end{pmatrix} \quad (2.11)$$

where now

$$\begin{aligned} \tilde{M}(x, E) &= \tilde{T}[x + (q-1)\beta, E] \dots \tilde{T}(x + \beta, E) \tilde{T}(x, E), \\ x_n &= 2\pi\beta q n + \delta = 2\pi p n + \delta. \end{aligned} \quad (2.12)$$

Now the transfer matrix \tilde{M} is independent of n . The eigenvalue condition on E is then just that the eigenvalues of the transfer matrix \tilde{M} lie on the unit circle. Since \tilde{T} , and therefore \tilde{M} , both satisfy

$$\det \tilde{T} = \det \tilde{M} = 1, \quad (2.13)$$

this condition becomes

$$2 \cos k = \text{Tr } \tilde{M}(\delta, E). \quad (2.14)$$

Having found the eigenvalues $E(\delta, k)$ using (2.14), the wavefunctions can be generated by means of the formula (2.9).

To summarise: there are two approaches to solving (2.1) when β is rational; one, which will be termed the Bloch picture, involves solving a $q \times q$ Hermitean eigenvalue equation, the other, which will be termed the Floquet picture, involves considering products of q 2×2 transfer matrices. The rest of this section will show how WKB methods can be applied when β is close to a rational number. First the application of the WKB method within the Bloch picture will be described. This has previously been attempted by Sokoloff (1981); it cannot be carried through to yield a full solution, but is worth describing since it is easier to understand because it is closer to ordinary WKB methods. Finally, the application of the WKB method in the Floquet picture will be described. This is harder to visualise, but does lead to a full solution of the problem.

When $\Delta\beta$ is small, the solution must 'locally' look like a solution of the form (2.3). On a 'global' scale, however, there is a slow change in the phase parameter δ ; in the region of the amplitude ψ_n the effective phase δ' is

$$\delta' = \delta + n\hbar/q, \quad (2.15)$$

where

$$\hbar = 2\pi\Delta\beta q. \quad (2.16)$$

The symbol \hbar is used in (2.15) because this quantity will be the small parameter of the WKB theory. In the neighbourhood of the amplitude ψ_n , the solution resembles a solution of the form (2.3) or (2.4) with δ replaced by δ' .

The Bloch wavevector, k , now varies slowly with n : the energy E is still given by equation (2.7), and is a constant for a given solution, so that (2.7) defines an implicit relationship between k and δ . The energy E should now be considered to depend on \hbar as well as δ and k ,

$$E = \varepsilon(\delta', k; \hbar) = \varepsilon_0(\delta', k) + \hbar\varepsilon_1(\delta', k). \quad (2.17)$$

since β in (2.7) depends on \hbar . The term of order \hbar in (2.17) will be important in what follows.

Following Sokoloff (1981), equation (2.1) is written in the form

$$\psi(x + \hbar/q) + \psi(x - \hbar/q) + 2\alpha \cos(2\pi px/\hbar + x - x_0)\psi(x) = E\psi(x), \quad (2.18)$$

where

$$x_0 = (2\pi p\delta/q) \bmod 2\pi, \quad \psi_n = \psi(x_n), \quad x_n = n\hbar/q + \delta. \quad (2.19)$$

By comparison with equation (2.5), this suggests a trial solution of the form

$$\psi(x) = A(x) \exp(iS(x)/\hbar) \sum_{m=0}^{q-1} a_m(x) \exp(2\pi ipmx/\hbar). \quad (2.20)$$

This trial solution corresponds to the abstract solution introduced in equation (1.1). The role of the vector \mathbf{u} in (1.1) is played by the set of Fourier coefficients a_m in (2.20). These coefficients are easily shown, by substituting (2.20) into (2.18), to satisfy the equation

$$\alpha e^{-ix} a_{m+1} + \alpha e^{ix} a_{m-1} + 2 \cos[(2\pi pm + S')/\hbar] a_m = Ea_m, \quad (2.21)$$

which is an eigenvalue equation for E corresponding to (1.2).

The next step in Sokoloff's approach to the WKB theory of Harper's equation is to expand $\psi(x \pm \hbar/q)$ in (2.20) in powers of \hbar , and insert the result into (2.18). Unfortunately, this does not lead to a consistent result; if the calculation is carried out correctly it is found that q independent equations are obtained which $A(x)$ should satisfy. (The solution which Sokoloff obtains for $A(x)$ is easily found to be incorrect.)

It turns out that a full solution of the WKB problem can be obtained in the Floquet picture, however. The transfer matrices $\tilde{M}(x, E)$ introduced in (2.11) are now no longer independent of n , but provided (2.2) is satisfied, these transfer matrices are at least slowly varying. It is possible to calculate the product of a string of slowly varying matrices;

$$\begin{aligned} \tilde{G}(x, x'; \hbar) &= \tilde{M}_E(x, \hbar) \tilde{M}_E(x - \hbar, \hbar) \dots \tilde{M}_E(x' + \hbar, \hbar) \tilde{M}_E(x', \hbar), \\ \tilde{M}_E(x, \hbar) &= \tilde{T}\{x + 2\pi p(q-1)/q + [(q-1)/q]\hbar, E\} \dots \tilde{T}(x + 2\pi p/q + \hbar/q, E) \tilde{T}(x, E), \\ \tilde{T}(x, E) &= \begin{pmatrix} E - 2\alpha \cos x & -1 \\ 1 & 0 \end{pmatrix}. \end{aligned} \quad (2.22)$$

A simple formula for the product $\tilde{G}(x, x'; \hbar)$ will be derived in § 3.

Before going on to discuss the WKB theory in detail, it will be helpful to describe briefly the final results of the calculation.

Suppose that β is a low denominator rational number, $\beta_0 = p/q$. The spectrum then consists of q bands (the centre two bands touch if q is even), and E is a periodic function of the Bloch wavevector k and position parameter δ ; $E = \varepsilon(\delta, k)$, with q branches, one for each band.

When β is close to β_0 , $\beta = \beta_0 + \hbar/2\pi q$, then WKB theory can be applied and δ and k become the position and momentum coordinates of the phase-space ($\delta \rightarrow x$, $qk \rightarrow p = S'$). The dispersion relation $E = \varepsilon(\delta, k)$ for a given band becomes the classical Hamiltonian; $H(x, p) = \varepsilon(x, kq)$.

When the phase trajectories (contours of $E = H(x, p)$) are closed orbits, the energies of the eigenstates are restricted by a Bohr-Sommerfeld quantisation condition. Some contours of a typical $H(x, p)$ are shown in figure 1 for the case $\alpha = 1$, when (by symmetry) all the phase trajectories are closed orbits. Each of the q bands of the

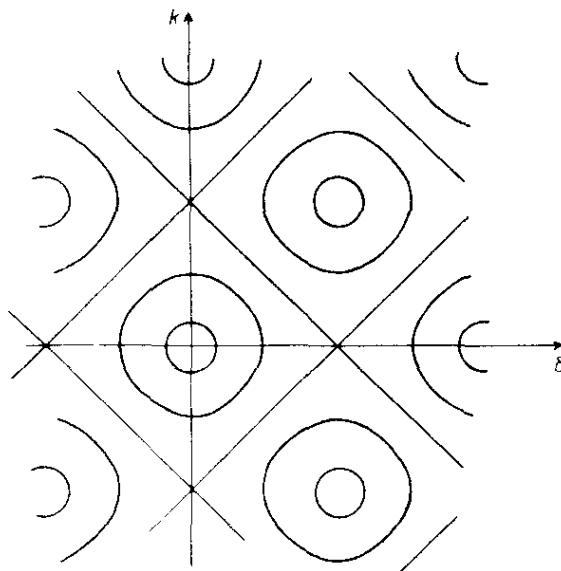
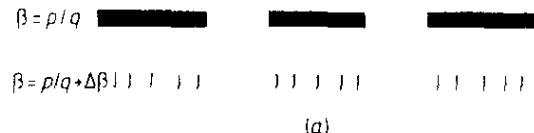
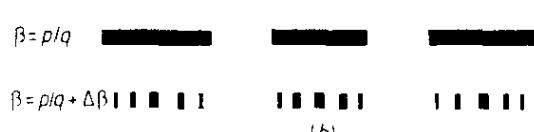


Figure 1. Phase trajectories of the classical Hamiltonian when $\alpha = 1$. All the phase trajectories are closed orbits



(a)



(b)

Figure 2. (a) Bloch bands of the commensurate system, $\beta = p/q$, become Bohr-Sommerfeld quantised levels when β is changed by a small amount. (b) The Bohr-Sommerfeld quantised levels are in practice broadened slightly by tunnelling effects.

spectrum is then split into a number of Bohr-Sommerfeld quantised levels. This situation is shown schematically in figure 2(a).

These Bohr-Sommerfeld quantised levels are not truly discrete; since the classical Hamiltonian $H(x, p)$ is a periodic function of x and p , the energy levels are broadened slightly by tunnelling between degenerate states, as illustrated schematically in figure 2(b). These tunnelling effects are discussed in detail in Wilkinson (1984) for the case $\beta_0 = 0$. Since the broadening of the levels due to tunnelling vanishes very rapidly as $\hbar \rightarrow 0$, as

$$E_{\text{tunnelling}} = \exp(-\text{constant}/\hbar), \quad (2.23)$$

it does still make sense to obtain an asymptotic formula for the Bohr-Sommerfeld quantisation condition.

3. Adiabatic matrix strings

This section derives a formula for the product of a string of slowly varying matrices, \tilde{M} . An asymptotic formula is obtained for the product \tilde{G} , defined by

$$\tilde{G}(x, x'; \hbar) = \prod_{\substack{n=0 \\ x'' = x' + nh}}^N \tilde{M}(x'', \hbar), \quad N = \frac{|x - x'|}{\hbar}, \quad (3.1)$$

i.e.

$$\tilde{G}(x, x'; \hbar) = \tilde{M}(x, \hbar) \tilde{M}(x - \hbar, \hbar) \dots \tilde{M}(x' + \hbar, \hbar) M(x', \hbar),$$

in the limit $\hbar \rightarrow 0$. Note that the matrices \tilde{M} depend on the slowness parameter of the adiabatic change, \hbar , as well as the variable x .

It is assumed that the matrices $\tilde{M}(x, \hbar)$ can be diagonalised

$$\tilde{M}(x, \hbar) = \tilde{X}^{-1}(x, \hbar) \tilde{D}(x, \hbar) \tilde{X}(x, \hbar), \quad (3.2)$$

where $\tilde{D}(x, \hbar)$ is diagonal, and that all the eigenvalues are distinct and lie on the unit circle.

Now (3.1) can be written

$$\tilde{G}(x, x'; \hbar) = \tilde{X}^{-1}(x, \hbar) \tilde{g}(x, x'; \hbar) \tilde{X}(x', \hbar) \quad (3.3)$$

where

$$\begin{aligned} \tilde{g}(x, x'; \hbar) &= \tilde{D}(x, \hbar) [\tilde{I} + \hbar \tilde{V}(x - \hbar, \hbar)] \tilde{D}(x - \hbar, \hbar) \\ &\times [\tilde{I} + \hbar \tilde{V}(x - 2\hbar, \hbar)] \dots [\tilde{I} + \hbar \tilde{V}(x', \hbar)] \tilde{D}(x', \hbar), \end{aligned} \quad (3.4)$$

and

$$\tilde{I} + \hbar \tilde{V}(x, \hbar) = \tilde{X}(x + \hbar, \hbar) \tilde{X}^{-1}(x, \hbar). \quad (3.5)$$

Throughout the calculation presented here, it will be sufficient to use the approximation

$$\begin{aligned} \tilde{V}(x) &= [d\tilde{X}(x, \hbar)/dx] \tilde{X}^{-1}(x, \hbar) + O(\hbar) \\ &= [d\tilde{X}(x, 0)/dx] \tilde{X}^{-1}(x, 0) + O(\hbar). \end{aligned} \quad (3.6)$$

(The second equation of (3.6) shows that, when calculating \tilde{V} , the dependence of \tilde{X} on \hbar can be neglected and will not be shown in subsequent equations.) Now, using the notation

$$\tilde{g}_0(x, x'; \hbar) = \tilde{D}(x, \hbar) \tilde{D}(x - \hbar, \hbar) \dots \tilde{D}(x' + \hbar, \hbar) \tilde{D}(x', \hbar), \quad (3.7)$$

and ordering the expansion of equation (3.4) in powers of \hbar

$$\tilde{g}(x, x'; \hbar) = \tilde{g}_0(x, x'; \hbar) + \hbar \sum_{n=0}^N \sum_{x''=x'+n\hbar} \tilde{g}_0(x, x'' + \hbar) \tilde{V}(x'') \tilde{g}_0(x'', x'; \hbar) + O(\hbar^2 V^2), \quad (3.8)$$

leads to an exact but implicit equation for \tilde{g}

$$\tilde{g}(x, x'; \hbar) = \tilde{g}_0(x, x'; \hbar) + \hbar \sum_{n=0}^N \sum_{x''=x'+n\hbar} \tilde{g}_0(x, x'' + \hbar, \hbar) \tilde{V}(x'') \tilde{g}(x, x'; \hbar). \quad (3.9)$$

An asymptotic solution of (3.9) will now be sought in the form

$$\tilde{g}(x, x'; \hbar) = \tilde{f}(x, x') \tilde{g}_0(x, x'; \hbar), \quad (3.10)$$

where $\tilde{f}(x, x')$ is diagonal. This trial solution is an adiabatic approximation; it expresses the expectation that when \hbar is small, so that \tilde{M} varies slowly, an eigenvector $u_i(x')$ of $M(x', \hbar)$ becomes, upon multiplying by $G(x, x'; \hbar)$, the corresponding eigenvector $u_i(x)$ of $M(x, \hbar)$. Before going any further, it is useful to define diagonal matrices \tilde{S} and \tilde{v} as follows

$$\tilde{D}(x, \hbar) = \exp[i\tilde{S}'(x, \hbar)], \quad (3.11)$$

$$\tilde{S}(x, x'; \hbar) = \int_{x'}^x dx'' \tilde{S}'(x'', \hbar),$$

$$\tilde{v}_j(x) = \begin{cases} \tilde{V}_j(x) & i=j \\ 0 & i \neq j. \end{cases} \quad (3.12)$$

Now, substituting (3.10) into (3.9), and making use of the definitions (3.11), (3.12)

$$\begin{aligned}
 g_{ij}(x, x'; \hbar) &= g_{0ij}(x, x'; \hbar) \\
 &= \hbar \sum_{x''} g_{0ij}(x, x'' + \hbar; \hbar) V_{ij}(x'') g_{0ij}(x'', x'; \hbar) f_j(x'', x') \\
 &= \hbar \sum_{x''} \exp\left(i \sum_{u=x''+\hbar}^x S'_i(u, \hbar)\right) V_{ij}(x'') f_j(x'', x') \exp\left(i \sum_{u=x'+\hbar}^{x''} S'_j(u, \hbar)\right) \\
 &= \hbar \exp[\frac{1}{2}i(S'_i(x) + S'_j(x'))] \sum_{x''} \exp[(i/\hbar)S_i(x, x'' + \hbar/2; \hbar)] V_{ij}(x'') f_j(x'', x') \\
 &\quad \times \exp[(i/\hbar)S_j(x'' + \hbar/2, x'; \hbar)] + O(\hbar) \\
 &= \exp[\frac{1}{2}i(S'_i(x) + S'_j(x'))] \times \int_{x'}^x dx'' V_{ij}(x'') f_j(x'', x') \\
 &\quad \times \exp[(i/\hbar)(S_i(x, x''; \hbar) + S_j(x'', x'; \hbar))] + O(\hbar). \tag{3.13}
 \end{aligned}$$

For terms with $i \neq j$, the integrand in (3.13) contains a rapidly oscillating term and gives a contribution of $O(\hbar)$, whereas for $i = j$ it gives a finite contribution. Therefore

$$\tilde{g}(x, x'; \hbar) = \tilde{g}_0(x, x'; \hbar) \left[\tilde{1} + \int_{x'}^x dx'' \tilde{f}(x'', x') \tilde{v}(x'') \right] + O(\hbar), \tag{3.14}$$

since only the diagonal elements of (3.13) remain. This justifies the use of the adiabatic approximation (3.10). From (3.10) and (3.14), \tilde{f} satisfies

$$\tilde{f}(x, x') = \tilde{1} + \int_{x'}^x dx'' \tilde{f}(x'', x') \tilde{v}(x''), \tag{3.15}$$

which gives

$$\tilde{f}(x, x') = \exp\left(\int_{x'}^x dx'' \tilde{v}(x'')\right), \tag{3.16}$$

where both \tilde{f} and \tilde{v} are diagonal. Therefore the central result of this section, the formula for $\tilde{G}(x, x'; \hbar)$, is found to be

$$\begin{aligned}
 \tilde{G}(x, x'; \hbar) &= \tilde{X}^{-1}(x) \exp\left(\int_{x'}^x dx'' \tilde{v}(x'')\right) \tilde{g}_0(x, x'; \hbar) \tilde{X}(x') + O(\hbar), \\
 \tilde{g}_0(x, x'; \hbar) &= \exp[\frac{1}{2}i(\tilde{S}'(x) + \tilde{S}'(x'))] \exp[(i/\hbar)S(x, x'; \hbar)] + O(\hbar).
 \end{aligned} \tag{3.17}$$

In this result the dependences of some quantities on \hbar have not been shown, since they are not important at this order of accuracy.

The remainder of this section will discuss a slight simplification of (3.17) which is possible when the transfer matrices preserve some quantity j , which will be called the current. This is usually the case in one-dimensional quantum mechanical problems. For any two vectors ϕ, ψ the current j is given by

$$j_{\phi, \psi} = \phi^{*\top} \tilde{J} \psi \tag{3.18}$$

(where \tilde{J} is a constant matrix). If j is preserved under the action of a transfer matrix \tilde{M} , then

$$j_{M\phi, M\psi} = (\tilde{M}\phi)^{\top} \tilde{J} (\tilde{M}\psi) = j_{\phi, \psi}, \tag{3.19}$$

so that M satisfies

$$\tilde{M}^{T*}\tilde{J}\tilde{M} = \tilde{J}. \quad (3.20)$$

The transfer matrices introduced in § 2 will be shown later to have this property. It can be shown that properties of \tilde{M} that are preserved on multiplying matrices together are indeed preserved by the formula (3.17) for the product \tilde{G} ; i.e. if it is real, unimodular, or satisfies the current conservation (3.20) then the approximate formula (3.17) for the product also has these properties.

It will be useful to get an expression for $\tilde{f}(x, x')$ in terms of the eigenvectors of $\tilde{M}(x, \hbar)$. Let $u_i(x)$ and $v_i(x)$ be right and left eigenvectors of $\tilde{M}(x, \hbar)$:

$$\tilde{M}u_i = \lambda_i u_i, \quad v_i \tilde{M} = \lambda_i v_i. \quad (3.21)$$

The u_i are proportional to the columns of \tilde{X}^{-1} and the v_i to the rows of \tilde{X} , so that

$$u_i \cdot v_j = N_{ij} \delta_{ij}. \quad (3.22)$$

For matrices \tilde{M} that satisfy (3.20), a useful relationship can be found connecting the left and right eigenvectors: from (3.20) it is easy to show that

$$(\tilde{J}u_i)^{*T}\tilde{M} = (\tilde{J}u_i)^{*T}\lambda_i^{*-1}, \quad (3.23)$$

so that the eigenvalues and eigenvectors come in pairs, related by

$$\lambda_i = \lambda_i^{*-1}, \quad v_i = (\tilde{J}u_i)^{*T}. \quad (3.24)$$

Since the eigenvalues λ_i are all on the unit circle for the transfer matrices considered here,

$$v_i = (\tilde{J}u_i)^{*T}. \quad (3.25)$$

Now collecting equations (3.16), (3.12), (3.6), (3.22), (3.25), a simple and useful formula can be given for $f_i(x, x')$, in terms of the eigenvector $u_i(x)$

$$\begin{aligned} f_i(x, x') &= \exp\left[\int_{x'}^x dx''(d\tilde{X}/dx|_{x''}\tilde{X}^{-1}(x''))_i\right] \\ &= \exp\left[-\int_{x'}^x dx''(u^{*T}\tilde{J} du/dx)/(u^{*T}\tilde{J} u)|_{x''}\right] \\ &= \exp\left[-\int_{x'}^x dx'' j_{u(du/dx)}/j_{uu}|_{x''}\right], \end{aligned} \quad (3.26)$$

where $u = u_i(x'')$ is the i th eigenvector of $M(x, \hbar)$.

4. Solution of the WKB problem for Harper's equation

This section uses the central results of § 3, equations (3.17) and (3.26), to solve the problem of finding a satisfactory WKB theory for Harper's equation (2.1).

The WKB solution required is of the form of equation (1.1)

$$\psi(x) = A(x)u(x) \exp\left(\frac{i}{\hbar} \int_{x'}^x p(x') dx'\right). \quad (4.1)$$

Locally, the solution can be described by a set of q amplitudes, either amplitudes of

a Bloch function u_n , or, equivalently, the Fourier components a_m of the Bloch function, as in equation (2.5). Alternatively, in the Floquet picture, ψ and u would be specified by just two amplitudes on a pair of adjacent lattice sites.

In the Floquet picture, by equation (3.17), if the wavefunction $\psi(x')$ is an eigenvector $u_i(x')$ at position x' , then at x it is given by

$$\begin{aligned}\psi(x) &= g_i(x, x'; \hbar) u_i(x) \\ &= f_i(x, x') g_{0i}(x, x'; \hbar) u_i(x).\end{aligned}\quad (4.2)$$

The $g_{0i}(x, x'; \hbar)$ term in (4.2), which can be written

$$g_{0i}(x, x'; \hbar) = \exp\left[\frac{i}{2}i(S'_i(x) + S'_i(x'))\right] \exp\left(\left(i/\hbar\right) \int_x^{x'} dx'' S'_i(x'', \hbar)\right), \quad (4.3)$$

can be associated with the oscillatory term in (4.1), and the $f_i(x, x')$ term with the amplitude $A(x)$.

The derivative $S'(x, \hbar)$ of the 'action' $S(x, x', \hbar)$ therefore plays the role of the momentum p of the phase space. The energy E is given as a function of $x (= \delta)$ and $S' (= kq)$ by (2.6) and (2.7) in the Bloch picture, or alternatively by (2.14) in the Floquet picture. The first order term in \hbar in the relations between E , x and S' must be retained since \hbar appears in the denominator of the argument of the exponential in (4.3).

It can be seen that the transfer matrix (2.10) satisfies the current conservation property (3.20), with

$$\tilde{J} = \begin{pmatrix} 0 & i \\ -i & 0 \end{pmatrix}. \quad (4.4)$$

The current $j_{\phi, \psi}$ can equivalently be calculated for the corresponding q dimensional vectors $|\phi\rangle, |\psi\rangle$

$$j_{\phi, \psi} = \langle \phi | \hat{J} | \psi \rangle, \quad (4.5)$$

where in the direct representation (by means of the Bloch function, u_n) the matrix elements of \tilde{J} are given by

$$J_{nn'} = \frac{i}{q} \begin{pmatrix} 0 & e^{-iS'/q} & & & & e^{iS'/q} & & \\ \cdot & \cdot & \cdot & \cdot & \cdot & 0 & & \\ & \cdot & \cdot & \cdot & \cdot & & \cdot & \\ & & \cdot & e^{iS'/q} & 0 & e^{-iS'/q} & & \\ & & & \cdot & \cdot & \cdot & \cdot & \\ 0 & & & & \cdot & \cdot & \cdot & \\ e^{-iS'/q} & & & & & \cdot & e^{iS'/q} & 0 \end{pmatrix}. \quad (4.6)$$

It is also useful to note that the eigenvalues and eigenvectors of the transfer matrices $M(x, E, \hbar)$ come in complex conjugate pairs, corresponding to points related by $\pm S'$ in the $x-S'$ plane. The eigenvalues and eigenvectors are therefore real when $S' = 0$:

$$u(x, S') = u^*(x, -S'). \quad (4.7)$$

If the eigenvectors u_i are given as functions of x , then $A(x)$ for the solution (4.2) is given by (cf (3.26))

$$A(x) = f_i(x, x') = \exp\left(-\int_{x'}^x dx'' j_{u(du/dx)/j_{uu}}\right), \quad (4.8)$$

and is in general complex. Alternatively $A(x)$ could be chosen to be real; this condition then defines a connection rule for the phases of the eigenvector $\mathbf{u}(x)$. If the $\mathbf{u}(x)$ are given with some arbitrary phase, then the correct phase is established by means of a transformation

$$\mathbf{u}(x) \rightarrow \mathbf{u}'(x) = \exp(i\phi(x))\mathbf{u}(x), \quad (4.9)$$

so that the current matrix element $j_{u(d\mathbf{u}/dx)}$ is transformed according to the equation

$$j_{u(d\mathbf{u}/dx)} \rightarrow j_{u'(d\mathbf{u}'/dx)} = j_{u(d\mathbf{u}/dx)} + i(d\phi/dx)j_{uu}. \quad (4.10)$$

Then by a suitable choice of $\phi(x)$, $j_{u(d\mathbf{u}/dx)}$ can be made to satisfy

$$\text{Im}(j_{u'(d\mathbf{u}'/dx)}) = 0, \quad (4.11)$$

so that by (4.8), $A(x)$ is now real.

5. The Bohr-Sommerfeld quantisation rule

This section derives the Bohr-Sommerfeld quantisation rule (1.3), which is the condition for single-valuedness of the WKB solution under continuation around a closed phase trajectory in the x, S' plane.

In this section it will be useful to consider the eigenvector $\mathbf{u}'(x, S')$ to be a given single-valued function defined on the phase plane. If $A(x)$ is to be a real function, the eigenvector \mathbf{u}' must be multiplied by a phase factor so that the modified eigenvector $\mathbf{u}(x, S')$ satisfies the connection formula (4.11), i.e.

$$\text{Im}[(j_{u(d\mathbf{u}/dx)})\Delta x + (j_{u(d\mathbf{u}/dS')})\Delta S'] = \text{Im}(j_{u\nabla u}) \cdot \Delta X = 0, \quad (5.1)$$

for transport of \mathbf{u} by a vector $\Delta X = (\Delta x, \Delta S')$ in the phase plane. This connection (5.1) is non-holonomic, and on transporting \mathbf{u} around a closed circuit C in phase space, it is multiplied by a phase factor $e^{i\gamma_c}$, given by

$$\begin{aligned} \gamma_c &= \text{Im} \oint_C (j_{u\nabla u}/j_{u'u'}) \cdot dX \\ &= \text{Im} \oint_C j_{u(d\mathbf{u}/dx)}/j_{u'u'} dx. \end{aligned} \quad (5.2)$$

On transporting \mathbf{u} around a phase trajectory of energy E , there is thus a phase change $e^{i\gamma(E)}$, where

$$\begin{aligned} \gamma(E) &= \text{Im} \oint_{E=E} (j_{u'\nabla u'}/j_{u'u'}) \cdot dX \\ &= \text{Im} \oint_{E=E} j_{u'(d\mathbf{u}'/dx)}/j_{u'u'} dx. \end{aligned} \quad (5.3)$$

This phase change makes a contribution to the Bohr-Sommerfeld quantisation formula.

In order to describe correctly the continuation of the solution around the phase trajectory of energy E , it is necessary to consider carefully what happens at the classical turning points, where $S' = 0 \bmod 2\pi$. On the lines $S' = 0$, the current j_{uu} is zero, since by (4.7), $\mathbf{u}(x, 0)$ is real. In the neighbourhood of the line $S' = 0$, $j_{uu}(x, S')$ takes the form

$$j_{uu} = a(x)S' + O(S'^2), \quad (5.4)$$

for some function $a(x)$. Now consider the form of $j_{u\nabla u}$ near the line $S'=0$. Consider the result

$$j_{u+\Delta u, u+\Delta u} = j_{uu} + 2 \operatorname{Re} j_{u\Delta u} + O(\Delta u^2) \quad (5.5)$$

(which uses the fact that the current operator is Hermitian), and take $u=u(x, 0)$ and $u+\Delta u=u(x+\Delta x, \Delta S')$. Then using (5.5) and ignoring $O(\Delta u^2)$

$$j_{u+\Delta u, u+\Delta u} = 2 \operatorname{Re} j_{u\Delta u} = 2 \operatorname{Re} j_{u\nabla u} \cdot \Delta X = a(x+\Delta x)\Delta S' \quad (5.6)$$

so that, near $S'=0$

$$\operatorname{Re}(j_{u\nabla u}) = (O, \frac{1}{2}a(x)) + O(x) + O(S'). \quad (5.7)$$

Now, from (4.8), the amplitude $A(x)$ (constrained to be real) is given by

$$A(x) = \exp\left(-\operatorname{Re} \int^x [j_{u(\Delta u/dx)}/j_{uu}]|_{x'} dx'\right). \quad (5.8)$$

Near the line $S'=0$, therefore,

$$A(x) \approx \exp\left[\int^x -\frac{1}{2}a(x')/(a(x')S'(x')) \cdot (dS'/dx') dx'\right] \quad (5.9)$$

$$A(x) \approx \text{constant}[S'(x)]^{-1/2}.$$

Thus $A(x)$ diverges at a classical turning point, x_0 , where $S'=0$. Near this point the form of the phase trajectory is given by

$$S'^2 = \text{constant}(x - x_0), \quad (5.10)$$

so that as x_0 is approached from within the classically allowed region, $A(x)$ diverges as

$$A(x) \sim \text{constant}(x - x_0)^{-1/4}. \quad (5.11)$$

(Of course there is not a real divergence of the exact solution, only in the WKB approximation; the assumption used in § 3 that the eigenvalues of the transfer matrix are distinct breaks down when $S'=0$.) The divergence of $A(x)$ given by (5.11) is of exactly the same type as is encountered in ordinary WKB problems at first-order turning points, and any of the usual arguments (e.g. continuation in the complex plane, see Landau and Lifshitz (1958)) show that an extra phase change of $\pi/2$ must be included for each of the two classical turning points of the phase trajectory.

The final contribution to the phase change of $\psi(x)$ is from the phase integral term: this is

$$\frac{1}{\hbar} \oint_{\epsilon(x, S')=E} S'(x, \hbar) dx. \quad (5.12)$$

As noted earlier, the correction to $\epsilon(x, s')$ of first order in \hbar must be retained when evaluating (5.12), since \hbar appears in the denominator. Collecting together all these contributions to the phase gives the Bohr-Sommerfeld quantisation rule for the system.

$$2\pi n = \frac{1}{\hbar} \oint_{\epsilon(x, S')=E} S'(x, \hbar) dx + \pi + \gamma(E). \quad (5.13)$$

Finally there are two important points which must be mentioned. Firstly, because j_{uu} is zero on the line $S'=0$, the integrand in the formula for $\gamma(E)$ diverges at the

classical turning point as $(x - x_0)^{-1/2}$. The integral $\gamma(E)$ remains finite, but does not tend to zero as E approaches a maximum or minimum of $\varepsilon(x, S')$, and the phase trajectory shrinks to a point. Instead, $\gamma(E)$ tends to a finite limit γ_0 at the top or bottom of a band. This limiting value of $\gamma(E)$ at the band edges is calculated in the appendix.

Secondly, there are some special cases which should be mentioned. When the rational number p/q to which β approximates is zero, then both the phase $\gamma(E)$ and the \hbar dependent corrections to $\varepsilon(x, S')$ vanish, and the Bohr-Sommerfeld quantisation condition takes the usual form. This case is discussed in detail in Wilkinson (1984). There are also some simplifications which occur when $p/q = \frac{1}{2}$, and it is only for p/q with denominators greater than two that all the effects described in this paper are seen.

6. Summary and comparison with numerical results

In this section some comparisons will be made of eigenvalues calculated using the Bohr-Sommerfeld quantisation rule with those calculated exactly. Firstly, however, the important formulae are collected together and summarised.

The parameter β in equation (2.1) is written

$$\beta = \beta_0 + \Delta\beta = p/q + \hbar/2\pi q. \quad (6.1)$$

The energy E is considered to be a function $\varepsilon(x, S')$ of the phase plane coordinates x and S' . The relationship between E , x and S' is, in the Bloch picture, given by the eigenvalue equation

$$\hat{H}(x, S')|u\rangle = E|u\rangle, \quad (6.2)$$

and the matrix elements of \hat{H} are given by (2.7), with $\delta \rightarrow x$ and $k \rightarrow S'/q$:

$H_{nn}(x, S')$

$$= \begin{pmatrix} 2\alpha \cos(x + 2\pi\beta) & e^{iS'/q} & & & e^{-iS'/q} \\ \cdot & \cdot & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ e^{-iS'/q} & 2\alpha \cos(x + 2\pi\beta n) & \cdot & \cdot & e^{iS'/q} \\ 0 & \cdot & \cdot & \cdot & \cdot \\ e^{iS'/q} & & e^{-iS'/q} & 2\alpha \cos(x + 2\pi\beta q) & \end{pmatrix}. \quad (6.3)$$

Equivalently, in the Floquet picture, this relationship is given by

$$2 \cos S' = \text{Tr } \tilde{M}_E(x, \hbar) \quad (6.4)$$

where

$$\tilde{M}_E(x, \hbar) = \tilde{T}(x + 2\pi\beta(q-1), E) \dots \tilde{T}(x + 2\pi\beta, E) \tilde{T}(x, E),$$

$$\tilde{T}(x, E) = \begin{pmatrix} E - 2\alpha \cos x & -1 \\ 1 & 0 \end{pmatrix}. \quad (6.5)$$

The phase change $\gamma(E)$ is given by a line integral in phase-space around a phase trajectory

$$\gamma(E) = \text{Im} \oint_{\varepsilon(x, S') = E} (j_u \nabla_u / j_{uu}) \cdot dX. \quad (6.6)$$

In the Bloch picture, the vector $\mathbf{u}(x, S')$ is an eigenvector of $\hat{H}(x, S')$, given by (6.3), and in the Floquet picture $\mathbf{u}(x, S')$ is an eigenvector of the transfer matrix $\hat{M}(x, \varepsilon(x, S'), \hbar)$. The matrix elements of the current operator are, in the Bloch picture

$$J_{nn} = \frac{i}{q} \begin{pmatrix} 0 & e^{-iS'/q} & & & e^{iS'/q} \\ \vdots & \ddots & \ddots & & \vdots \\ & \ddots & \ddots & 0 & e^{-iS'/q} \\ & & & 0 & \ddots \\ e^{-iS'/q} & & & & e^{iS'/q} & 0 \end{pmatrix}, \quad (6.7)$$

and in the Floquet picture

$$\tilde{J} = \begin{pmatrix} 0 & i \\ -i & 0 \end{pmatrix}. \quad (6.8)$$

The final result, the Bohr-Sommerfeld quantisation condition, is (cf 5.13)

$$\oint_{E=F} S'(x, E, \hbar) dx = 2\pi[n + \frac{1}{2} - (1/2\pi) \operatorname{sign}(\hbar) \gamma(E)] \cdot |\hbar| \quad (6.9)$$

(remember that \hbar can be negative for this system).

Now the theoretical predictions of this paper will be compared with some numerical results.

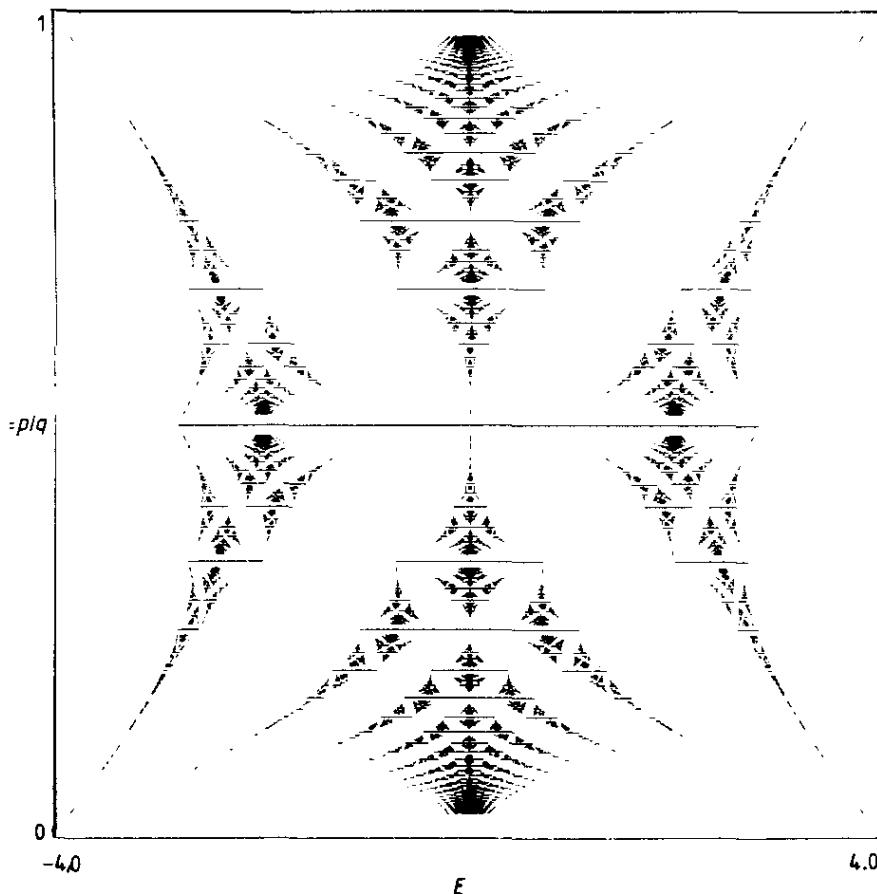


Figure 3. A plot of the spectrum of Harper's equation, plotted for every rational $\beta = p/q$ with $q \leq 40$. This picture illustrates the situation shown schematically in figure 2.

First, figure 3 gives a general illustration of the ideas discussed in this paper. This picture, originally published by Hofstadter (1976), is a plot of the spectrum (i.e. the support of $\varepsilon(\delta, k)$) of Harper's equation for every rational value of β with denominator q less than 40. There are q energy bands (the central two touching when q is even). When β is close to a low-denominator rational, p_0/q_0 , then the q bands are very narrow and cluster into q_0 groups, which correspond to the q_0 bands when $\beta = \beta_0 = p_0/q_0$. These q narrow bands are the Bohr-Sommerfeld levels discussed in this paper, slightly broadened by tunnelling effects.

It is also easy to see by inspection of figure 3 that, except for $\beta_0 = 0, 1$ and $\beta_0 = \frac{1}{2}$, the Bohr-Sommerfeld quantisation is not of the usual form, i.e. action $= 2\pi(n + 1/2)\hbar$, since the pattern is not symmetrical about and below the lines $\beta = \beta_0$ (cf. equation (6.9), which is not symmetric under $\hbar \rightarrow -\hbar$).

Next, table 1 gives some comparisons of energy levels predicted using formulae (6.1)–(6.9), E_{pr} , with exact levels, E_{ex} obtained by the same method as Hofstadter, using high precision arithmetic. In all cases except one the broadening of these levels by tunnelling is smaller than the six-digit precision of the eigenvalues E_{ex} . The meaning of the columns p_0, q_0 , is given by equation (6.1), n is the quantum number of the Bohr quantised level (as in equation (6.9)) and n_b labels the energy band within which this level lies. The subscripts t and b of n_b mean that n is counted from the top and bottom of the energy band respectively. Finally, ΔE is the separation of the closest neighbouring energy level, and gives a scale against which the error $|E_{ex} - E_{pr}|$ of the predictions should be compared.

Finally, table 2 shows how the error in the predicted levels decreases very rapidly as \hbar decreases. These results suggest that the error of the prediction is $O(\hbar^3)$ compared to the separation of levels which decreases as $O(\hbar)$.

Table 1. Comparison of predicted energy levels E_{pr} with exact levels E_{ex} . For full description see § 6 of text. For all values in this table, $\alpha = 1$.

p_0	q_0	$\Delta\beta$	n_b	n	E_{pr}	E_{ex}	$ E_{pr} - E_{ex} $	ΔE	$\gamma(E)$
1	3	1/200	1t	1	-1.99435	-1.99277	0.00158	0.08395	1.6791
1	3	1/200	1t	6	-2.29301	-2.29209	0.00092	0.04154	0.6868
1	3	1/200	1t	10	-2.40938	-2.40860	0.00078	0.02295	-0.4313
1	3	1/200	2b	1	-0.61889	-0.62003	0.00114	0.10129	-1.6406
1	3	1/200	3t	1	2.70954	2.71090	0.00136	0.03379	-2.2567
1	3	1/200	3t	5	2.58201	2.58442	0.00241	0.02839	-3.5757
1	4	1/300	2t	1	-0.28742	-0.28171	0.00621	0.10782	1.4364
1	4	1/300	2t	5	-0.58727	-0.58578	0.00148	0.04299	0.0375
1	3	2/387	1t	1	-1.99414	-1.99252	0.00162	0.08631	1.6786
1	3	2/387	1t	7	-2.33540	-2.33427 to -2.33432	0.001	0.035	0.4193
3	7	1/1960	3b	1	-1.58070	-1.58068	0.00002	0.00623	-0.5570
3	7	1/1960	3b	5	-1.60464	-1.60460	0.00004	0.00557	-0.6380

7. Concluding remarks

This paper has demonstrated a novel type of Bohr-Sommerfeld quantisation, involving a non-holonomic connection rule for transporting the eigenvectors u of the matrix-valued Hamiltonian function $\hat{H}(x, p)$ around a circuit in phase space.

Table 2. Illustrating the rapid improvement of the predictions as \hbar decreases. All the results in this table refer to the case $\alpha = 1$, $p_0 = 1$, $q_0 = 3$, $n_b = 3t$. Using equation (A9), the limiting value of $\gamma(E)$ at the edge of the band concerned is predicted to be $\gamma_0 = -1.8200$. The values of $\gamma(E)$ for the $n = 1$ states approach this limiting value.

n	$\Delta\beta$	E_{pr}	E_{ex}	$ E_{pr} - E_{ex} $	ΔE	$\gamma(E)$
1	1/100	2.68295	2.69043	0.00748	0.06637	-2.9737
	1/200	2.70954	2.71090	0.00136	0.03379	-2.2567
	1/400	2.72120	2.72140	0.00020	0.01710	-1.9760
	1/800	2.72671	2.72671	0.00000	0.00859	-1.8554
5	1/200	2.58201	2.58442	0.00241	0.02839	-3.5757
	1/400	2.65491	2.65514	0.00023	0.01575	-2.3862
	1/800	2.69286	2.69288	0.00002	0.00826	-2.0428

The phase change γ of the eigenvector has been determined in terms of the line integral of the connection (5.2). In principle, γ could also be expressed as the integral of the curvature of the connection over the area enclosed by C . In practice, however, this is not useful, since this curvature is singular on the line $S' = 0$, and in any case for a computer calculation of $\gamma(E)$ the line integral is much easier to evaluate.

The method given in this paper is easily adapted to the problem of determining the Bohr-Sommerfeld quantisation condition for Bloch electrons in a weak magnetic field. It is well known that this condition takes the form

$$\mathcal{A}_k = 2\pi(eB/\hbar)(n + \Gamma), \quad (7.1)$$

where \mathcal{A}_k is the area enclosed by a section through the Fermi surface perpendicular to the magnetic field (Onsager 1952). The constant Γ is not determined by Onsager's argument, and has previously only been determined exactly by very elaborate methods based on the effective Hamiltonian approach (see e.g. Roth 1966). The main result is that, for crystals with centres of inversion, Γ is always equal to $\frac{1}{2}$ (plus terms of higher order in the magnetic field), but when there is not a centre of inversion there is an additional component of Γ given by an integral analogous to (5.3).

It is worthwhile to note that two other results have appeared recently which involve a non-holonomic connection rule for the phase of the eigenvector of a matrix.

Firstly Mead and Truhlar (1979) and Berry (1984) calculate the phase change of the wavefunction of a system after being varied slowly around a cyclic path in the space of some parameters of the system. In terms of the matrix product calculation of § 3 of this paper, this corresponds to considering a string of slowly varying unitary evolution operators, i.e. to the case $\tilde{J} = \tilde{I}$ in equation (3.20).

Secondly Thouless *et al* (1982) have considered the quantised Hall effect in samples with a weak periodic potential with a rational number of flux quanta per unit cell. They show that the Hall conductance of a full sub-band is $e^2/2\pi\hbar$ times the phase change when the wavefunction is transported around the edge of the magnetic Brillouin zone using the adiabatic connection rule of Mead and Berry. Because the magnetic Brillouin zone is topologically a torus, they are able to show that this phase change is 2π times an integer. Simon (1983) has exhibited a connection between the work of Berry and Mead, and that of Thouless *et al*, and has emphasised the importance of the idea of a non-holonomic connection.

The formula given in § 3 for the product of a string of slowly varying matrices is a very general result, and may have many uses other than those considered here.

Acknowledgment

I wish to thank the UK Science and Engineering Research Council for a postgraduate studentship.

Appendix

This appendix demonstrates the result that $\gamma(E)$ approaches a finite limit γ_0 at the top or bottom of an energy band, by calculating this limit in terms of the transfer matrices.

The transfer matrix is assumed to be known in the form of a series expansion:

$$\tilde{M}(x, E, \hbar) = \begin{pmatrix} A & B \\ C & D \end{pmatrix} = \begin{pmatrix} A_0 & B_0 \\ C_0 & D_0 \end{pmatrix} + (x - x_0) \begin{pmatrix} a_x & b_x \\ c_x & d_x \end{pmatrix} + (E - E_0) \begin{pmatrix} a_E & b_E \\ c_E & d_E \end{pmatrix} + \hbar \begin{pmatrix} a_\hbar & b_\hbar \\ c_\hbar & d_\hbar \end{pmatrix} + \dots, \quad (\text{A1})$$

where x_0 and E_0 are the values of x and E corresponding to the top or bottom of the energy band $E = \epsilon(x, s')$ when $\hbar = 0$. Also, it is assumed that the first few coefficients in the expansion of $\text{Tr } \tilde{M}$ are known

$$2 \cos S' = \text{Tr } \tilde{M}(x, E, \hbar) = 2 + r(x - x_0)^2 - s(E - E_0) - t\hbar. \quad (\text{A2})$$

From equation (2.8), it can be seen that $r = q^2 \alpha^q$. From (A2), the momentum S' is given by

$$S'^2 = s(E - E_0) - r(x - x_0)^2 - t\hbar, \quad (\text{A3})$$

in the neighbourhood of the band edge. The action $S(E)$ for the phase trajectory of energy E is given by

$$S(E) = \oint_{\epsilon = E} S' dx = (\pi/\alpha^q q)[s(E - E_0) - t\hbar]. \quad (\text{A4})$$

Now the phase change $\gamma(E)$ of the eigenvector \mathbf{u} of \tilde{M} will be found, in the neighbourhood of the energy E_0 . The transfer matrix \tilde{M} of (A1) can be diagonalised as follows

$$\tilde{M}(x, E, \hbar) = \tilde{X}^{-1} \tilde{D} \tilde{X}, \quad \tilde{D} = \begin{pmatrix} e^{iS'} & 0 \\ 0 & e^{-iS'} \end{pmatrix} \quad (\text{A5})$$

$$\tilde{X} = \begin{pmatrix} C & D - e^{-iS'} \\ -C & e^{iS'} - D \end{pmatrix}, \quad \tilde{X}^{-1} = \frac{1}{2iC \sin S'} \begin{pmatrix} e^{iS'} - D & e^{-iS'} - D \\ C & C \end{pmatrix}.$$

Now, for the eigenvector corresponding to the eigenvalue $e^{iS'}$, the differential element of the phase change γ is (retaining only lowest order terms)

$$d\gamma_1 = \text{Im } d\tilde{V}_{11} = \text{Im}(d\tilde{X} \tilde{X}^{-1})_{11} = \frac{D-1}{2CS'} dC - \frac{dD}{2S'}. \quad (\text{A6})$$

Using the relationship (A3) between x , S' and E , for any quantity $A = A(x, E)$

$$\left(\frac{\partial A}{\partial S'} \right)_x = \frac{2S'}{s} \left(\frac{\partial A}{\partial E} \right)_x, \quad \left(\frac{\partial A}{\partial x} \right)_{S'} = \left(\frac{\partial A}{\partial x} \right)_E + \frac{2\alpha^q q^2 (x - x_0)}{s} \left(\frac{\partial A}{\partial E} \right)_x. \quad (\text{A7})$$

Using (A7) in (A6), to lowest order

$$\begin{aligned} d\gamma_1 &= \gamma_x dx + \gamma_{S'} dS', \\ \gamma_x &= \{[(D_0 - 1)/2C_0]c_x - \frac{1}{2}d_x\}/S' + O(1) = \kappa/S' + O(1), \\ \gamma_{S'} &= O(1). \end{aligned} \quad (\text{A8})$$

Also, it is easy to show that $d\gamma_2(x, -S') = d\gamma_1(x, S')$. The phase change is then given by, in the limits $x, S', \hbar \rightarrow 0$:

$$\gamma(E) \approx \gamma_0 = \oint_{\epsilon=E} \gamma_x dx + \gamma_{S'} dS' \approx \oint \kappa/S' dx \quad (\text{A9})$$

i.e.

$$\gamma_0 = 2\pi\kappa/\alpha^{q/2}q, \quad \kappa = [(D_0 - 1)/2C_0]c_x - \frac{1}{2}d_x. \quad (\text{A10})$$

The Bohr-Sommerfeld quantisation condition in the neighbourhood of the band edge then becomes (using (5.13), (A4), (A9))

$$E = E_0 + (1/s) \cdot 2\alpha^q q \hbar [n + \frac{1}{2} - (\kappa/\alpha^{q/2}q) \operatorname{sign}(\hbar) + t/2\pi], \quad (\text{A11})$$

where the level number satisfies the inequality

$$n > (\kappa/\alpha^{q/2}q) \cdot \operatorname{sign}(\hbar) - \frac{1}{2}. \quad (\text{A12})$$

References

- Bellissard J and Simon B 1982 *J. Funct. Anal.* **48** 408–19
 Berry M V 1984 *Proc. R. Soc. A* **392** 45–57
 Harper P G 1955 *Proc. Phys. Soc. A* **68** 874–8
 Hofstadter D R 1976 *Phys. Rev. B* **14** 2239–49
 Khinchin A Ya 1964 *Continued Fractions* (Chicago: University Press)
 Landau L D and Lifshitz E M 1958 *Quantum Mechanics* ch 7 (Oxford: Pergamon)
 Mead C A 1979 *J. Chem. Phys.* **70** 2276–83
 Mead C A and Truhlar D G 1979 *J. Chem. Phys.* **70** 2284–96
 Onsager L 1952 *Phil. Mag.* **43** 1006–8
 Roth L M 1966 *Phys. Rev.* **145** 434–48
 Simon B 1982 *Adv. Appl. Maths.* **3** 463–90
 —— 1983 *Phys. Rev. Lett.* **51** 2167–70
 Sokoloff J B 1981 *Phys. Rev. B* **23** 2039–41
 Thouless D J, Kohmoto M, Nightingale M P and den Nijs M 1982 *Phys. Rev. Lett.* **49** 405–9
 Wilkinson M 1984 *Proc. R. Soc. A* **391** 305–50

Current address: Department of Physics and Applied Physics, John Anderson Building, University of Strathclyde, Glasgow G4 0NG, Scotland, United Kingdom.

Chapter 7

WEss-ZUMINO TERMS AND ANOMALIES

- [7.1] M. Stone, "Born-Oppenheimer Approximation and the Origin of Wess-Zumino Terms: Some Quantum-Mechanical Examples," *Phys. Rev.* **D33** (1986) 1191 361
- [7.2] J. Goldstone and F. Wilczek, "Fractional Quantum Numbers on Solitons," *Phys. Rev. Lett.* **47** (1981) 986 365
- [7.3] E. Witten, "Global Aspects of Current Algebra," *Nucl. Phys.* **B223** (1983) 422 369
- [7.4] I. J. R. Aitchison, "Berry Phases, Magnetic Monopoles, and Wess-Zumino Terms or How the Skyrmeion got its Spin," *Acta Phys. Polonica* **B18** (1987) 207–235 380
- [7.5] P. Nelson and L. Alvarez-Gaumé, "Hamiltonian Interpretation of Anomalies," *Commun. Math. Phys.* **99** (1985) 103–114 409



7

Wess–Zumino Terms and Anomalies

In this section we look at two topics with origins in quantum field theory, Wess–Zumino terms and anomalies, from the special point of view afforded by geometric phases. We trust that readers who have followed us this far will find such an approach especially accessible. The fundamental ideas have close analogies with phenomena we have already met in previous chapters, and have found applications well outside their original scope, so that we hope that readers with other backgrounds will not be scared away. Most of the papers we have selected for this section have explicitly pedagogical intentions. Also, we have tried to choose a graded cross-section, so that a determined reader can reach the heights by manageable steps. Since the papers do speak so well for themselves, in this introduction we shall merely define a few of the special terms and make some general observations. Those who wish to delve further may wish to consult Refs. [1] and [2].

The original Wess–Zumino term was invented, by Wess and Zumino,³ in constructing an effective action for mesons. Its purpose was to describe certain interactions, such as $\pi^0 \rightarrow \gamma\gamma$, that the naive effective action failed to include. Since then, it has been generalized to a number of other situations. Without writing a Wess–Zumino term down explicitly (see [7.2] for details) we can summarize some of its most important properties:

- (i) The Wess–Zumino term is first order in time derivatives, and persists in the adiabatic limit.
- (ii) Its coefficient is quantized, for topological reasons. This is a generalization of the famous Dirac quantization condition, which shows that the existence of magnetic monopoles is only consistent with the principles of quantum mechanics, if the magnetic charges of all particles are integer multiples of a certain basic charge.
- (iii) It plays a crucial role in determining the kinematics of quantization. It directly modifies the canonical momentum, and, furthermore, the value of its coefficient modulo 2 dictates whether certain solitons of the model, which represent baryons, are to be quantized as bosons or fermions.

- (iv) It results from “integrating out” very high-energy or high-mass states (quarks, in this case). These states may leave behind unusual quantum numbers (the elementary solitons possess baryon number 1) for collective states of the effective low-energy theory.
- (v) It has a singular, gauge non-invariant form when written down explicitly as a term in the effective Lagrangian. However, it may also be expressed in a manifestly invariant and non-singular form by adding an extra spatial dimension. That is, one considers a five-dimensional manifold \mathcal{M} whose boundary is four-dimensional (compact Euclidean) spacetime, and writes the corresponding piece of the action as the integral over \mathcal{M} of a well-behaved total divergence.

Reader, you will recognize that any one of these properties could also refer to the term we added to the effective Born–Oppenheimer Lagrangian in the introduction to Chapter 4 and [4.3], in order to describe the effect of Berry’s phase on nuclei in diatomic molecules. We saw there that, because the typical splitting between different electronic configurations is very large in comparison with nuclear splittings, we can validly integrate out the electrons in the Born–Oppenheimer approximation, but there remains a highly non-trivial residue of these lost degrees of freedom. The relative phases of the electron ground-state configurations for different nuclear positions must be taken into account when sewing them together to form wavefunctions for the nuclear degrees of freedom. The net effect of the phases is to generate an extra term in the effective nuclear action

$$S_{\text{phase}} = \int A(R) \cdot \dot{R} dt = \int A \cdot dR \quad (7.1)$$

involving a magnetic monopole gauge potential. This term satisfies properties (i) and (iv). It satisfies (ii) by calculation (ultimately because of the quantization condition of electronic angular momentum), and (iii) because, as we saw in Chapter 4, it alters the quantization of the nuclear angular momentum. Finally, because of the Dirac string singularity, the term (7.1) in the Lagrangian is not gauge invariant, although the integrated action along a closed loop is invariant up to the addition of a multiple of 2π when the Dirac quantization condition is satisfied. One way to see this is by using Stokes’ theorem to write

$$S_{\text{phase}} = e \iint_S dR_i dR_j F_{ij} \quad (7.2)$$

where S is a surface whose boundary is the original path in parameter space and $F_{ij} = g\epsilon_{ijk}r^k/|r|^3$ is the gauge-invariant magnetic field of a monopole. The phase is now manifestly gauge-invariant, but seems to depend on the bounding surface. However, the Dirac condition that the product of electric

and magnetic charge eg be a multiple of $1/2$, implies that the ambiguity is a multiple of 2π , and does not affect the phase $e^{iS_{\text{phase}}}$.

Given that the term (7.1) satisfies all of the properties enumerated above, perhaps we should refer to it as a particular example of a Wess-Zumino term. Viewed in this light, the WZ term is nothing more than a field-theoretic generalization of a very basic feature of the Born-Oppenheimer approximation. This point of view is discussed further, in light of explicit examples, in the enclosed paper by Stone [7.1]. Aitchison's paper [7.2] gives a more leisurely and detailed account of the connection between Berry's phase and Wess-Zumino terms, and includes several explicit examples.

One often employs a type of Born-Oppenheimer approximation in the quantization of collective states—solitons—in quantum field theory. Consider, for instance, a theory in which light scalar fields interact with much heavier fermions. One expects that as long as the scalar fields vary slowly in space-time, there will be a large gap between the ground state of the fermion system and any other state of this system, and furthermore that this ground state may be constructed by appropriately patching together locally determined ground-states. The situation is altogether analogous to the one we just met, in connection with the quantization of diatoms. In the paper included here, Goldstone and Wilczek [7.3] analyze the field theory problem. They show that in this situation the scalar field configurations do generally inherit quantum numbers from the fermions. Although they did not use this language, the charges and currents they compute can be summarized by the effective Lagrangian one would obtain for the scalar fields by integrating out the fermions. This effective Lagrangian would contain a term of Wess-Zumino type, associating fermion quantum numbers to certain configurations, and instructing us to quantize these as fermions.

It was in the context of showing how a Lagrangian containing only (bosonic) meson degrees of freedom—which on general grounds is expected to be the effective theory for QCD at low energies—could also contain (fermionic) baryon degrees of freedom as collective excitations, that Witten wrote the remarkable and very influential paper reprinted here [7.4].

Next let us say a few words concerning anomalies. A symmetry is said to be anomalous if it exists at the level of the classical theory, but is inevitably spoiled upon quantization. A particularly well-studied class of anomalies occurs in quantizing theories containing fermions, and is closely related to Wess-Zumino terms as discussed above. One fruitful way of looking at these anomalies is very reminiscent of the Born-Oppenheimer approximation, if we think of the fermions as “fast” and the boson fields as “slow” degrees of freedom. This separation need have nothing to do with the relative mass scales of the two types of fields; here, it is really just a convenience that allows us to construct the theory in stages, by first defining the fermion states

for fixed configurations of the other fields, performing the integration over fermionic variables, and then quantizing the resulting effective Lagrangian for the bosonic fields. In particular, we need to define the Dirac sea associated with each Bose field configuration, and patch these together. Now it is often found that holonomy, of the kind that by now should seem quite familiar, can arise on the configuration space. To maintain classical symmetries of the system, we would like to relate (parallel transport) the Dirac seas for different configurations by making symmetry transformations. It is often found, however, that a non-integrable phase holonomy prevents us from doing this consistently. It can also occur that the transport law demanded by one symmetry conflicts with that demanded by another, so that at most one of them can be maintained on quantization.

Some elementary but illuminating remarks introducing anomalies are contained in the paper by Jackiw [1.2]. The paper by Alvarez-Gaumé and Nelson [7.5] is a lovely but very sophisticated treatment from the point of view of Berry's phase. We recommend that you read all the other papers in this section before attempting this one; after that preparation we think you will find it accessible and rewarding.

Anomalies can be good: as will be discussed below, some have even been observed experimentally. On the other hand, anomalies in local gauge symmetries are widely believed to be unacceptable—a belief we share. For local gauge invariance is necessary to insure that longitudinal photons, and their relative gluons of other kinds, are not present in the physical spectrum. And they'd better not be, because they are ghosts—the probability of creating one, if not zero, is negative! So cancellation of anomalies in all local gauge symmetries is usually imposed as a constraint on possible quantum field theories. However, anomalous *global* symmetries do not make a field theory inconsistent; they simply get broken by the quantization process. When fermions associated with a global anomaly are integrated out, the effective action for the remaining degrees of freedom will generally include a Wess-Zumino term. Its role is to account for the above-mentioned holonomy.

In conclusion, it may be appropriate to say a few words about the physical applications of anomalies. They find several important applications in QCD, the modern theory of the strong interaction. First of all there is the original application, to the description of the decay $\pi^0 \rightarrow \gamma\gamma$. A naive application of chiral symmetry, which had been spectacularly successful in describing much of the low-energy phenomenology of pions, led to the very poor prediction of a vanishing rate for this decay.⁴ It was therefore tremendously important when Adler, and Bell and Jackiw⁵ demonstrated that an anomaly spoiled chiral symmetry in this decay. Moreover, Adler⁶ traced this anomaly down to a single graph, and showed that the decay rate was profoundly and quantitatively related to the fundamental field content of whatever theory

described the strong interaction. The successful quantitative prediction of the $\pi^0 \rightarrow 2\gamma$ rate, which is sensitive to the color and fractional charge assignments of quarks, still stands as one of the most remarkable triumphs of quantum field theory. Recently, the related process $\gamma\gamma \rightarrow \pi\gamma$ has also been measured, and found to agree with the anomaly-based prediction.

QCD with massless quarks, which is a good approximation for most hadronic physics, is formally scale invariant. However this classical symmetry is spoiled upon quantization. The “scale anomaly” is parametrized by an effective coupling constant that varies depending upon the characteristic energy and momentum scales of the process under consideration.⁷ Of course, it was the observed fact that scale invariance is approximately but not exactly realized among strongly interacting particles that led to the discovery of asymptotic freedom and QCD in the first place,⁸ and to many successful quantitative predictions for high-energy processes.

QCD with massless quarks also seems to contain another symmetry, axial baryon number, that is not observed in nature. Furthermore, there is no sign of its being spontaneously broken (that is, there is no light Nambu-Goldstone boson with the appropriate quantum numbers).⁹ It was therefore most welcome when t’Hooft¹⁰ demonstrated that this apparent symmetry is in fact anomalous, that is not a symmetry the quantized theory at all. With this discovery, the match between the symmetries predicted by QCD and the observed symmetries of the strong interaction in nature became perfect, and the last significant barrier to acceptance of QCD as the theory of the strong interaction fell.

Because anomalies are often sensitive to arbitrarily massive degrees of freedom, and because they depend only on very general features of the theory involved, they have played an important role in guiding speculations concerning physics beyond the standard model. The cancellation of anomalies in electroweak gauge symmetries occurs in a non-trivial manner between quarks and leptons. This was one of the first suggestions of deep connections among these particles, and inspired efforts toward unification. It also goes at least part of the way toward explaining why quarks and leptons occur in repeating families, each with the same basic structure.

Finally, much of the recent interest in string theory was stimulated by the discovery of Green and Schwarz¹¹ that certain anomalies which spoiled the consistency of the theory cancelled for special choices of space-time dimension and internal gauge group (10 and SO(32), respectively). Their mechanism for cancellation of anomalies also worked at a formal level for the exceptional group $E_8 \times E_8$; however at the time no way of incorporating this gauge group into a string theory was known. Within a few weeks, inspired by this cancellation, Gross, Harvey, Martinec and Rohm¹² discovered the heterotic string construction, which does incorporate $E_8 \times E_8$. Today, most attempts to connect string theory with observed reality start from the

heterotic string.

- [1] S. Treiman, R. Jackiw, B. Zumino, and E. Witten, *Current Algebra and Anomalies* (Princeton: University Press, 1985).
- [2] A.P. Balachandran, "Wess-Zumino Terms and Quantum Symmetries," and A. Niemi, "Quantum Holonomy," in M. Jezabek and M. Praszałowicz, eds., *Skyrmions and Anomalies* (Singapore: World Scientific, 1987); A. Dhar and S. Wadia, *Phys. Rev. Lett.* **52** (1984) 959.
- [3] J. Wess and B. Zumino, *Phys. Lett.* **36B** (1971) 95.
- [4] D. G Sutherland, *Nucl. Phys.* **B2** (1967) 433.
- [5] S. Adler, *Phys. Rev.* **177** (1969) 2426.
J. S. Bell and R. Jackiw, *Nuovo Cimento* **60A** (1969) 47.
- [6] S. Adler, in *Lectures on Elementary Particles and Quantum Field Theory*, Proc. 1970 Brandeis Summer Institute, ed. S. Deser *et al.* (Cambridge: MIT Press, 1970).
- [7] D. Gross, in *Methods in Field Theory*, Proc. of 1975 Les Houches Summer School, ed. R. Balian and J. Zinn-Justin (Amsterdam: North-Holland, 1976).
- [8] D. Gross and F. Wilczek, *Phys. Rev. Lett.* **30** (1973) 1343;
H. Politzer, *Phys. Rev. Lett.* **30** (1973) 1346.
- [9] S. Weinberg, *Phys. Rev.* **D11** (1975) 3583.
- [10] G. 't Hooft, *Phys. Rev. Lett.* **59B** (1976) 172.
- [11] M. Green and J. Schwarz, *Phys. Lett.* **149B** (1984) 285.
- [12] D.J. Gross, J. Harvey, E. Martinec, and R. Rohm, *Nucl. Phys.* **B256** (1985) 253–284.

Born-Oppenheimer approximation and the origin of Wess-Zumino terms: Some quantum-mechanical examples

Michael Stone

*Department of Physics, University of Illinois at Urbana-Champaign,
1110 West Green Street, Urbana, Illinois 61801*

(Received 30 September 1985)

I provide some simple quantum-mechanical examples in which the Berry phase gives rise to Wess-Zumino terms. The connection with the Born-Oppenheimer approximation is also discussed.

I. INTRODUCTION

When, in a path-integral description of a quantum field theory, we have a system of Fermi fields interacting with bosons our first reaction is to eliminate the Grassmann variables needed for the fermions, integrating them out to produce the Mathews-Salam determinant.¹ This determinant is a functional of the Bose fields and acts back on them in the form of an effective action. Often, especially when the fermions are heavy, the effective action just renormalizes the parameters of the Bose action,² but sometimes the results are more startling: Symmetries may be lost and the theory may even be inconsistent. These unexpected effects are due to terms in the effective action of the type discussed by Wess and Zumino.³ They encapsulate all the anomalies of the system. A great deal of progress has recently been made in understanding the topological character and origins of these terms.^{4,5} In particular, Nelson and Alvarez-Gaume⁶ have explained their origin in a Hamiltonian context.

It seems worthwhile to show just how simple the explanation in Ref. 6 is by giving a few examples from quantum mechanics. In this simple context one can calculate everything yet there is still enough structure to retain the essentials of their ideas.

In Sec. II, I shall exhibit some nontrivial fiber bundles that originate naturally in extremely simple models. In Sec. III, I will quantize them by a path-integral route, which mimics what we do in a field theory. Then we repeat the quantization in the Born-Oppenheimer approximation where we will note some deficiencies in the accounts of this method given in textbooks. A conventional solution will be presented in the Appendix so that the reader can verify that nothing inconsistent with the usual rules of quantum mechanics is being perpetrated.

II. MONOPOLES FROM FERMIONS

Consider the family of Hamiltonians $H(\hat{n})$, labeled by points on S^2 (i.e., \hat{n} is a three-dimensional unit vector),

$$H(\hat{n}) = \mu \hat{n} \cdot \sigma . \quad (2.1)$$

Here σ are the usual Pauli matrices. For each \hat{n} there is a two-dimensional Hilbert space. Let us select a basis vector in each one corresponding to the eigenvalue $+\mu$. We can obtain these eigenstates by the use of spin-projection opera-

tors

$$|\psi_+^{(1)}\rangle = N^{\frac{1}{2}}(1 + \hat{n} \cdot \sigma) \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} \cos\theta/2 \\ \sin(\theta/2)e^{i\phi} \end{pmatrix} , \quad (2.2)$$

$$|\psi_+^{(2)}\rangle = N^{\frac{1}{2}}(1 + \hat{n} \cdot \sigma) \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} \cos(\theta/2)e^{-i\phi} \\ \sin\theta/2 \end{pmatrix} ,$$

where N is a real normalization factor and $\hat{n} = (\sin\theta \cos\phi, \sin\theta \sin\phi, \cos\theta)$. Clearly,

$$|\psi_+^{(1)}(\hat{n})\rangle = e^{i\phi} |\psi_+^{(2)}(\hat{n})\rangle . \quad (2.3)$$

We need (at least) two choices of $|\psi_+\rangle$ because $|\psi_+^{(1)}\rangle$ is a reasonable eigenvector everywhere except at the south pole (where the phase is ambiguous) and $|\psi_+^{(2)}\rangle$ is reasonable everywhere except at the north pole. No one choice of phase is smooth everywhere on S^2 ; the one-dimensional rays at each \hat{n} fit together to form a nontrivial bundle⁷ with (2.3) as the transition function between two patches (if we had not insisted on choosing subspaces the total Hilbert space bundle is trivial). We can characterize the eigenspace bundle by putting a connection on it which is compatible with (2.3). One such connection which will arise naturally in Sec. III is

$$iA_{\perp}^{(1)} = \langle \psi_+^{(1)}(\hat{n}) | d\psi_+^{(1)}(\hat{n}) \rangle = \frac{i}{2}(-\cos\theta + 1)d\phi , \quad (2.4a)$$

$$iA_{\perp}^{(2)} = \langle \psi_+^{(2)}(\hat{n}) | d\psi_+^{(2)}(\hat{n}) \rangle = \frac{i}{2}(-\cos\theta - 1)d\phi . \quad (2.4b)$$

$A_{\perp}^{(1)}$ and $A_{\perp}^{(2)}$ differ by a gauge transformation,

$$A_{\perp}^{(1)} = A_{\perp}^{(2)} + d\phi . \quad (2.5)$$

so both have the same curvature

$$F_{\perp} = dA_{\perp} = \frac{1}{2}\sin\theta d\theta \wedge d\phi = \frac{1}{2}d(\text{area}) \quad (2.6)$$

(the A_{\perp}, F_{\perp} corresponding to eigenvalue $-\mu$ just differ in sign from these). We recognize the gauge potentials and fields of the monopole bundle⁷ with $\int F = 2\pi$.

Another example comes from restricting \hat{n} to lie in the (x, z) plane so $H(\hat{n})$ is real. The eigenvectors can be chosen real,

$$\psi_+ = \begin{pmatrix} \cos\theta/2 \\ \sin\theta/2 \end{pmatrix} , \quad \psi_- = \begin{pmatrix} \sin\theta/2 \\ -\cos\theta/2 \end{pmatrix} , \quad (2.7)$$

so A is identically zero—but we have to use the transition

function

$$|\psi_{\pm}(2\pi)\rangle = -|\psi_{\pm}(0)\rangle$$

to make the choice single valued. The bundle is still twisted, being the Möbius strip.

The twisted nature of these bundles is only a curiosity so far; it will become important when we raise \hat{n} to the status of a dynamical variable in Sec. III.

III. A SPINNING SOLENOID

We will study the Hamiltonian

$$H = \frac{1}{2I} L^2 - \mu \sigma \cdot \hat{n}, \quad (3.1)$$

where L is the angular momentum operator that generates rotations of \hat{n} . One could imagine constructing a Heath-Robinson-type device which would be described by this Hamiltonian: A long thin solenoid rotating about its center of mass would have $\frac{1}{2}L^2/I$ as its Hamiltonian and placing a spin- $\frac{1}{2}$ particle with magnetic moment μ at this center of mass would produce the second term. When μ is small the two systems would spin independently. As μ becomes large the spin will become slaved to the direction of the solenoid and its spin- $\frac{1}{2}$ nature will be transferred to the orbital dynamics of the solenoid. It is the large- μ case which we will consider here because I wish to describe the induction of Wess-Zumino terms by the failure of the naive decoupling theorem⁸ and large μ corresponds to a large fermion mass.

A path integral for (3.1) is

$$\int d[\hat{n}] d[\bar{\psi}] d[\psi] \delta(n^2 - 1) \times \exp \left[- \int \left(\frac{I\hbar^2}{2} + \bar{\psi}(\partial_t - \mu \hat{n} \cdot \sigma) \psi \right) dt \right]. \quad (3.2)$$

The Fermi operators allow four states: no spin, one spin up, one spin down, and two spins—one of each, up and down. The ground state for large μ will be one spin aligned against the field. We will evaluate the fermion determinant by noting that large μ enables us to use the adiabatic theorem since the ground state is well separated from any other state to which transitions may be stimulated by the motion of \hat{n} .

The adiabatic theorem⁹ says that for slowly varying $H(t)$ we can approximate the solution of the time-dependent Schrödinger equation

$$i\partial_t |\psi\rangle = H(t) |\psi\rangle \quad (3.3)$$

in terms of eigenstates $|\psi^0\rangle$ of the "snapshot" Hamiltonian

$$H(t) |\psi^0(t)\rangle = E(t) |\psi^0(t)\rangle \quad (3.4)$$

as

$$|\psi(t)\rangle = \exp \left\{ -i \int_0^t E(t') dt' + i\gamma(t) \right\} |\psi^0(t)\rangle. \quad (3.5)$$

where γ is the "Berry phase"^{9,10} obeying

$$i\frac{d\gamma}{dt} + \langle \psi^0 | \frac{d}{dt} \psi^0 \rangle = 0. \quad (3.6)$$

The determinant for a closed path Γ of duration β is just

the vacuum-vacuum amplitude so

$$\det[\partial_t + H(n(t))] = \exp \left[- \int_0^\beta E(t) dt + i\gamma(\Gamma) \right]. \quad (3.7)$$

We have changed to Euclidean time in conformity with (3.2). (Note γ , as a phase, does not change.) The effect of the spin is to change the \hat{n} path integral to

$$\int d[\hat{n}] \delta(n^2 - 1) \exp \left[- \int \left(\frac{I\hbar^2}{2} + E(t) + \langle \psi(\hat{n}) | \frac{\partial \psi}{\partial t}(\hat{n}) \rangle \right) dt \right] \quad (3.8)$$

$$= \text{const} \times \int d[n] \delta(n^2 - 1) \exp \left[- \int \left(\frac{I\hbar^2}{2} + iA_+(\hat{n}) \cdot \dot{n} \right) dt \right]. \quad (3.9)$$

The dynamics of \hat{n} has become the motion of a charged particle on a sphere, moving under the influence of a magnetic monopole at the center. The $A \cdot \dot{n}$ in (3.9) is the simple example of a "Wess-Zumino" term introduced by Witten in Ref. 4. It is not rotationally invariant and cannot be written globally without introducing string singularities although its dynamical effects are both rotationally invariant and non-singular. In this example, unlike Ref. 4, the Wess-Zumino term has been induced dynamically instead of being inserted by hand and is a consequence of the phase ambiguities and bundle structure of Sec. II.

This result is sufficiently remarkable as to require a separate derivation (or two). One can obtain the same physics without the path integral by using the Born-Oppenheimer approximation. This says that if we have a system of slowly moving "nuclear" coordinates R , and some fast moving "electronic" coordinates r whose Hamiltonian $H_e(R)$ is parametrized by R , then the wave function $\Phi(r, R)$ can again be expressed in terms of solutions to the "snapshot" Hamiltonian

$$H_e(R) \psi(r, R) = E(R) \psi(r, R) \quad (3.10)$$

as

$$\Phi(r, R) = \phi(R) \psi(r, R), \quad (3.11)$$

and ϕ obeys the modified Schrödinger equation

$$-\frac{\hbar^2}{2M} \nabla^2 \phi + [E(R) + V(R)] \phi(R) = i\hbar \frac{\partial}{\partial t} \phi(R) \quad (3.12)$$

where ∇ is a covariant derivative:

$$\nabla_\mu = \frac{\partial}{\partial R^\mu} + \left\langle \psi | \frac{\partial \psi}{\partial R^\mu} \right\rangle = \partial_\mu + iA_\mu. \quad (3.13)$$

This approximation has exactly the same physics as the adiabatic theorem and is made by ignoring the same sets of matrix elements. It is worth noting that textbooks usually omit the gauge potentials. They have been discussed in the literature,¹¹ however, and are crucial here.

This factorization of the total wave function shows up the source of the gauge dependence of $\phi(R)$. The total wave function is gauge independent—but a change of phase of $\psi(r, R)$ must be compensated for by a change of phase in $\phi(R)$ which thus inherits the phase problems from the integrated-out fermions. The path integral, with the fermions replaced by the effective action, is the path integral

for $\phi(R)$, nor for the total wave function.

In our case $\phi(R)$ is the wave function for free motion on a sphere around the monopole, i.e., the monopole harmonics, so the eigenstates of (3.1) are

$$\langle \hat{n}, s | \phi \rangle = [D_{m=1/2}^l(\theta, \phi, \psi)]^* \langle s | \psi_+(\hat{n}) \rangle , \quad (3.14)$$

$$E_{jm} = \frac{1}{2I} j(j+1) + \text{const} , \quad (3.15)$$

degeneracy = $2j+1$.

Just as the ordinary Y_m^l are equal to $[D_{mo}^l(\theta, \phi, 0)]^*$ (the D 's being the ordinary rotation matrices in Euler angle parametrization), the monopole harmonics for a particle of charge (q times unit charge) in orbit about a monopole are $[D_{m=1/2}^l(\theta, \phi, \psi)]^*$. The dependence on the angle ψ reflects the gauge dependence. One must choose an angle ψ for each θ, ϕ and one cannot make this choice [of a point in $SU(2) = S^3$] under the inverse of the Hopf map: $S^3 \rightarrow S^2$ without introducing gauge patches. That a constrained spin is both classically and quantum-mechanically equivalent to motion about a monopole has been discussed by Leinaas in Ref. 12.

The second example, restricting the motion to the (x, z) plane, can be dealt with similarly. It is essentially the same as the problem of quantizing the collective modes of the $SU(2)$ Skyrmion. There is no Wess-Zumino term since $A = 0$ but the homotopy group of the space of paths is $\pi_1(S_1) = Z$ and we have to choose phases for the different classes of configurations. The quantum mechanics clearly requires $(-1)^n$, where n is the element of $\pi_1(S_1)$ involved. In the Born-Oppenheimer approach the states are

$$\langle \theta, s | j \rangle = e^{i\theta} \langle s | \psi_+(\theta) \rangle , \quad j = n + \frac{1}{2} . \quad (3.16)$$

Each factor is double valued but the product is not.

IV. DISCUSSION

These quantum-mechanical models provide a simple realization of how attempting to decouple a degree of freedom, by making the energy gap too large to excite it, can fail and effects are left behind. This is what happens in a quantum field theory when one makes fermion masses large in the

way discussed in Ref. 8. These models are also useful for thinking about the non-Abelian [or $\pi_2(G_3)$] and global [or $\pi_1(G_3)$] anomalies in the spirit of Ref. 6 (which was the main influence behind this paper). The analogy can be stretched too far, however: The analysis presented here is only valid in the adiabatic limit. If the coupling between the spin and the solenoid is such that the spin can be excited, the excited state has a Berry phase with the opposite sign and the topological effects are washed out (in other words, our total Hilbert space is a trivial product of the spin and configuration spaces and it is only in the adiabatic limit that we are forced into a subspace with nontrivial twist). This does not happen in a quantum field theory which is much more subtle—at least in the cases of interest. The Berry phase in the field theory arises as a property of the Dirac sea; the ground state and all excitations built on it have the same Berry phase, provided that the vacuum structure is not completely disrupted, and so the Wess-Zumino terms do not depend on the adiabatic limit which is then relegated to being a convenient trick to compute them.¹³

ACKNOWLEDGMENTS

This work was supported by the National Science Foundation Grant No. DMR84-15063. I would also like to thank NORDITA for hospitality and my colleagues at Urbana for encouraging me to give the talks on Wess-Zumino terms which led to my seeking a pedagogical example that even I could understand.

APPENDIX: CONVENTIONAL APPROACH TO QUANTUM MECHANICS

The problem of the solenoid coupled to a spin can, of course, be solved in a conventional manner. I shall do this here.

The Hamiltonian is

$$H = \frac{1}{2I} \mathbf{L}^2 - \mu \hat{\mathbf{n}} \cdot \boldsymbol{\sigma} \quad (A1)$$

and is clearly invariant under combined rotations of $\boldsymbol{\sigma}$ and $\hat{\mathbf{n}}$ so $[H, J] = 0$, $J = \mathbf{L} + \frac{1}{2}\boldsymbol{\sigma}$. This suggests using the states

$$\begin{aligned} |j - \frac{1}{2}, j, m\rangle &= \frac{1}{\sqrt{2j}} [(j+m)^{1/2} |j - \frac{1}{2}, m - \frac{1}{2}\rangle | \frac{1}{2}\rangle + (j-m)^{1/2} |j - \frac{1}{2}, m + \frac{1}{2}\rangle | - \frac{1}{2}\rangle] , \\ |j + \frac{1}{2}, j, m\rangle &= \frac{1}{\sqrt{2j+2}} [-(j-m+1)^{1/2} |j + \frac{1}{2}, m - \frac{1}{2}\rangle | \frac{1}{2}\rangle + (j+m+1)^{1/2} |j + \frac{1}{2}, m + \frac{1}{2}\rangle | - \frac{1}{2}\rangle] . \end{aligned} \quad (A2)$$

which are eigenstates of L^2 , J^2 , and J_z expressed in terms of eigenstates of L^2 , L_z , and S_z .

The operator $(\boldsymbol{\sigma} \cdot \hat{\mathbf{n}})$ is odd under $\hat{\mathbf{n}} \rightarrow -\hat{\mathbf{n}}$, commutes with J , and obeys $(\boldsymbol{\sigma} \cdot \hat{\mathbf{n}})^2 = 1$. So

$$\boldsymbol{\sigma} \cdot \hat{\mathbf{n}} |j \pm \frac{1}{2}, j, m\rangle = -|j \mp \frac{1}{2}, j, m\rangle . \quad (A3)$$

(The minus sign requires some computation since the properties of $\boldsymbol{\sigma} \cdot \hat{\mathbf{n}}$ quoted above could have given ± 1 .) Restricted to the space with fixed j, m , H becomes

$$H_{jm} = \frac{1}{2I} [j(j+1) + \frac{1}{4}] + \begin{pmatrix} (j + \frac{1}{2})/2I & \mu \\ \mu & -(j + \frac{1}{2})/2I \end{pmatrix} \quad (A4)$$

so

$$E_{jm} = \frac{1}{2I} [j(j+1) + \frac{1}{4}] \pm \left[\frac{(j + 1/2)^2}{4I^2} + \mu^2 \right]^{1/2} . \quad (A5)$$

For large μ the ground state is

$$|E_{jm}^-, j, m\rangle = \frac{1}{\sqrt{2}} (|j + \frac{1}{2}, j, m\rangle - |j - \frac{1}{2}, j, m\rangle) . \quad (A6)$$

We can express this in the factorized form of (3.14) by applying some rotation operators: Write

$$|E_{jm}^-, j, m\rangle = \begin{pmatrix} \phi_{+1/2}(\hat{\mathbf{n}}) \\ \phi_{-1/2}(\hat{\mathbf{n}}) \end{pmatrix} .$$

where $\phi_s(\hat{n}) = (\hat{n}, s | E_{jm}^-, j, m)$. The action of a rotation operator $U(R)$ is

$$\begin{aligned} U(R)|\hat{n}, s\rangle &= |R\hat{n}, s'\rangle D_{ss'}^{1/2}(R) , \\ U(R)|E_{jm}^-, j, m\rangle &= |E_{jm}^-, jm'\rangle D_{mm'}^j(R) , \end{aligned} \quad (A7)$$

so the matrix element $(\hat{n}, s | U(R) | E_{jm}^-, j, m)$ can be evaluated in two ways to give

$$\begin{aligned} (\hat{n}, s | E_{jm}^-, j, m') D_{mm'}^j(R) &= D_{ss'}^{1/2}(R^{-1}\hat{n}, s' | E_{jm}^-, jm') \\ (\hat{n}, s | E_{jm}^-, j, m) &= D_{ss'}^{1/2}(R)(R^{-1}\hat{n}, s' | E_{jm}^-, jm') D_{mm'}^j(R^{-1}) . \end{aligned} \quad (A8)$$

Arrange $R^{-1}\hat{n}$ to be the north pole and use the properties

of the angular momentum states,

$$(\theta, \phi | J \pm \frac{1}{2}, m \pm \frac{1}{2}) = Y_m^{J \pm \frac{1}{2}}(\theta, \phi)$$

at the north pole

$$Y_m^l(\theta = 0, \phi) = \delta_{m,0} \left[\frac{2l+1}{4\pi} \right]^{1/2}$$

to see that

$$(R^{-1}\hat{n}, s' | E_{jm}^-, j, m) \propto \delta_{s'm} \delta_{m,1/2} .$$

Using this we can write

$$\begin{aligned} (\hat{n}, s | E_{jm}^-, j, m) &= \langle s | \psi_+(R) \rangle D_{1/2,m}^j(R^{-1}) \\ &= [D_{m,1/2}^j(R)]^* \langle s | \psi_+(R) \rangle \end{aligned} \quad (A9)$$

as promised by Eq. (3.14). Different choices of R , the rotation which takes the north pole to \hat{n} , give rise to a different distribution of the phases between the two factors.

¹P. T. Mathews and A. Salam, Nuovo Cimento **12**, 563 (1954), 2, 120 (1955).

²T. Appelquist and J. Carazzone, Phys. Rev. D **11**, 2856 (1975).

³J. Wess and B. Zumino, Phys. Lett. **37B**, 95 (1971).

⁴E. Witten, Nucl. Phys. **B223**, 422 (1983).

⁵There have been many papers on this topic. See B. Zumino, Nucl. Phys. **B253**, 477 (1985), and references therein.

⁶P. Nelson and L. Alvarez-Gaume, Commun. Math. Phys. **99**, 103 (1985).

⁷One does not need to know the language of fiber bundles to understand this paper but it helps to put things in context. A good review is T. Eguchi, P. Gilkey, and A. J. Hanson, Phys.

Rep. **66**, 213 (1980).

⁸E. d'Hoker and E. Farhi, Nucl. Phys. **B248**, 59 (1984), **B248**, 77 (1984).

⁹See L. I. Schiff, *Quantum Mechanics*, 3rd ed. (McGraw-Hill, New York, 1968), p. 289.

¹⁰M. Berry, Proc. R. Soc. London **A392**, 45 (1984), B. Simon, Phys. Rev. Lett. **51**, 2167 (1983).

¹¹C. A. Mead and D. G. Truhlar, J. Chem. Phys. **70**, 2284 (1979).

¹²J. M. Leinaas, Phys. Scr. **17**, 483 (1978).

¹³Similar points are made in C. Gomez, Phys. Rev. D **32**, 2235 (1985).

Fractional Quantum Numbers on Solitons

Jeffrey Goldstone^(a)

Stanford Linear Accelerator Center, Stanford University, Stanford, California 94305

and

Frank Wilczek

Institute for Theoretical Physics, University of California, Santa Barbara, California 93106

(Received 9 July 1981)

A method is proposed to calculate quantum numbers on solitons in quantum field theory. The method is checked on previously known examples and, in a special model, by other methods. It is found, for example, that the fermion number on kinks in one dimension or on magnetic monopoles in three dimensions is, in general, a transcendental function of the coupling constant of the theories.

PACS numbers: 11.10.Lm, 11.10.Np

Peculiar quantum numbers have been found to be associated with solitons in several contexts:
 (i) The soliton provides, of course, a different background than the usual vacuum around which to quantize other fields. The difference between these "vacuum polarizations" may induce unusual quantum numbers localized on the soliton.¹⁻³
 (ii) Solitons may require unusual boundary conditions on the fields interacting with them, in particular leading to conversion of internal quantum numbers into rotational quantum numbers.¹⁻⁶

(iii) In the case of dyons, there is classically a family of solitons with arbitrary electric charge. The determination of which of these are in the physical spectrum requires quantum-mechanical considerations and brings in the θ parameter of non-Abelian gauge theories.^{7,8}

At present all these phenomena seem distinct although there are suggestive relationships. In this note, we shall concentrate on (i), proposing a general method of analysis and working out a few examples.

An intuitively appealing, and perhaps physically realizable, example of the phenomena we are addressing are the fractionally charged solitons on polyacetylene.^{2,3,9} A caricature model of a polyacetylene molecule is shown in Fig. 1(a)—in the ground state we have alternating single and double bonds, which may be arranged in two inequivalent but degenerate forms *A* and *B*. If there is an imperfection, as shown in Fig. 1(b), we go from *A* on the left-hand side to *B* on the right-hand side. This configuration cannot be brought to either pure *A* or pure *B* by any finite rearrangement of electrons, and so it will relax to a stable configuration—a soliton. If we put two imperfections together, as in Fig. 1(c), we find a configuration which begins and ends as *A*. Compared to the corresponding segment of pure *A*, it is missing one bond. If we add an electron to the two-imperfection strand, we can deform this configuration by a finite rearrangement into a pure *A* strand. (We are pretending, for simplicity, that each bond represents a single electron instead of a pair.) Interpreting this, we see that a two-soliton state is equivalent to the ground state if we add an electron. Thus, by symmetry, each separated soliton must carry electron number $-\frac{1}{2}$ (and electric charge $+\frac{1}{2}e$).

We can relate these stick-figure pictures of polyacetylene to field theory as follows: Let $d_1 > d_2$ be the internuclear distances characterizing single and double bonds, respectively. Define a scalar field which is a function of the link *i* by $\varphi_i = (-1)^i(d - \frac{1}{2}d_1 - \frac{1}{2}d_2)$, where *d* is the internuclear distance for link *i*. Thus in the *A* configuration $\varphi_i = \frac{1}{2}(d_1 - d_2)$ (independent of *i*), in the *B* configuration $\varphi_i = -\frac{1}{2}(d_1 - d_2)$, and in the soliton configuration φ_i interpolates between these values. Now we can show that it makes sense to approximate φ_i by a continuum field and the interactions of the electrons with φ (a charge-density wave) by $\mathcal{L}_I = g\bar{\psi}\gamma^5\varphi\psi$; furthermore the electrons

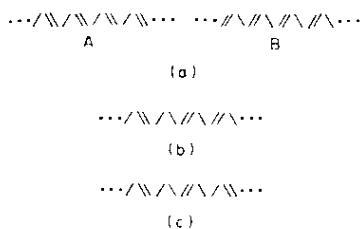


FIG. 1. (a) The two degenerate ground states for electronic structure of polyacetylene. (b) An imperfection interpolating between the two ground states. (c) A chain with two imperfections.

can be treated for present purposes (near the Fermi energy) as relativistic particles.

In this formulation, we make contact with the work of Jackiw and Rebbi.¹ They found that the spectrum of the Dirac equation in the presence of a soliton contains a zero-energy solution. By symmetry, this solution is composed of (projects onto) half a positive-energy and half a negative-energy solution with respect to the normal ground state. Thus if we fill the zero-energy level, we have a soliton state with electron number $+\frac{1}{2}$; if we leave it empty, the electron number is $-\frac{1}{2}$.

Su and Schrieffer have described a generalization,¹⁰ which occurs in a chain with a repeating unit of single-single-double bonds, as in Fig. 2. A slight modification of the discussion of Fig. 1 shows that we now have solitons which can be added in triples to give the normal ground state, deficient by one electron. We expect the electron number of a single soliton to be $-\frac{1}{3}$.

A field theoretic model must now have essentially new features. Jackiw and Rebbi emphasized that in their model the Dirac equation in the presence of a soliton has a charge-conjugation symmetry, and then their interpretation of the zero modes cannot account for any charges other than half-integral. Thus we will consider models where the background destroys all symmetries which interchange positive- and negative-energy solutions of the Dirac equation.

Our method of calculating the soliton quantum number will be to imagine building up the soliton by slow changes in fields, starting from the ground state. In order to reach the solitons by slow changes, we may have to enlarge the field space during intermediate stages, as we shall see. In any case, for slow variations of fields in space and time, we can readily compute the flow of the appropriate charge in the no-particle state. We then simply integrate to find the accumulated charge on the soliton.

Let us illustrate these remarks on a concrete example. We consider, in $1+1$ dimensions, massless fermions interacting with two scalar fields φ_1 and φ_2 as follows:

$$\mathcal{L}_I = g\bar{\psi}(\varphi_1 + i\gamma_5\varphi_2)\psi. \quad (1)$$

Now if φ_1 and φ_2 are slowly varying in space and

...

FIG. 2. A form of polymer with single-single-double bond pattern in the ground state.

time, i.e., their gradients are $\ll g(\varphi_1^2 + \varphi_2^2)^{1/2}$, we may conveniently calculate the change in the expected value of $j^\mu = \bar{\psi} \gamma^\mu \psi$ in the no-particle state by considering the Feynman graph of Fig. 3. Since the interaction (1) is chirally invariant, we may first suppose that only $\varphi_1 \neq 0$ at a given point, and then express the result in a chirally symmetric form. We then need only do a very simple calculation for an effectively massive fermion to find

$$\begin{aligned} \langle j^\mu \rangle &= \frac{1}{2\pi} \epsilon^{\mu\nu} \epsilon_{ab} \frac{\varphi_a \partial_\nu \varphi_b}{|\varphi|^2} \\ &= \frac{1}{2\pi} \epsilon^{\mu\nu} \partial_\nu \tan^{-1} \frac{\varphi_2}{\varphi_1}, \end{aligned} \quad (2)$$

If the scalar fields do not propagate (they represent very massive particles) more complicated graphs need not be considered.

If in the end we reach the soliton state by slow changes, we need only to evaluate (2) to find the fermion number charge on the soliton. It is important to remark that the resulting state will be a true eigenstate of the charge, not a superposition of states of different charge (even though we only derived an expectation value). For this it is only necessary to note that there are no degenerate states of different charge. In this the localized charge on a soliton differs from, for instance, the "localized charges" of $\frac{1}{2}$ on the top and bottom of an ammonia ion.

Two general features of the result deserve comment. First, the divergence $\partial_\mu j^\mu$ vanishes identically, reflecting the conservation of fermion number. Second, the charge $Q = \int d^3x \langle j^0 \rangle = (2\pi)^{-1} \times \Delta(\tan^{-1} \varphi_2/\varphi_1)$ is independent of the coupling constant g and depends only on the values of φ_1 and φ_2 at spatial infinity.

We can represent a massive fermion by fixing $\varphi_1 = m/g$. If the theory supports a soliton for which $\varphi_2(x) \rightarrow \pm i$ as $x \rightarrow \pm\infty$, we find

$$Q = \pi^{-1} \tan^{-1}(gr/m). \quad (3)$$

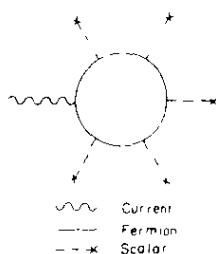


FIG. 3. Vacuum polarization graphs for evaluation of induced currents.

Notice that this is a transcendental function of the couplings! As $m \rightarrow 0$, we find $Q = \frac{1}{2}$; this is the Jackiw-Rebbi case of a single (linear) scalar coupling. The limit $m \rightarrow 0$ is delicate just because there are two degenerate states of charge $\pm \frac{1}{2}$ in the limit. If we take $m = 0$ from the beginning, adiabatic changes will fill these equally on an average. The current would vanish. A slight perturbation lifts the degeneracy. Of course the charge $-\frac{1}{2}$ state is reached by letting $m \rightarrow 0$ through negative values.

A field theory version of the chains of Figs. 1 and 2 is the interaction

$$\mathcal{L}_I = g \bar{\psi} e^{i\theta \gamma_5} \psi \quad (4)$$

for which we find

$$\langle j^\mu \rangle = (2\pi)^{-1} \epsilon^{\mu\nu} \partial_\nu \theta, \quad Q = (2\pi)^{-1} \Delta \theta. \quad (5)$$

The solitons with θ varying from 0 to π (so two together give $0 \rightarrow 2\pi \sim 0$, equivalent to vacuum) have charge $\frac{1}{2}$; with θ varying from 0 to $2\pi/3$, charge $\frac{1}{3}$; etc.

Some 1+1 dimensional models become especially transparent if the method of bosonization is employed. In 1+1 dimensions, one can rewrite fermion fields as nonlocal expressions in boson fields.¹¹ Some bilinears transform in a simple local way, however:

$$\begin{aligned} i\bar{\psi} \gamma^\mu \partial_\mu \psi &\rightarrow \frac{1}{2} \partial_\mu \varphi \partial^\mu \psi, \\ \bar{\psi} \gamma^\mu \psi &\rightarrow \epsilon^{\mu\nu} \partial_\nu \psi / \sqrt{\pi}, \\ \bar{\psi} \psi &\rightarrow \mu \cos 2\sqrt{\pi} \varphi, \\ i\bar{\psi} \gamma_5 \psi &\rightarrow \mu \sin 2\sqrt{\pi} \varphi \end{aligned}$$

(μ is an arbitrary scale parameter). Thus the interaction (4) becomes in this representation $\mathcal{L}_I = g \mu \cos(2\sqrt{\pi} \varphi - \theta)$. Now if θ in a soliton varies by $\Delta\theta$ from $-\infty$ to $+\infty$, the potential $-\mathcal{L}$ is minimized when $\varphi = \theta/2\sqrt{\pi}$; in particular, $\Delta\varphi = \Delta\theta/2\sqrt{\pi}$. Integrating $\bar{\psi} \gamma^0 \psi = \partial_0 \psi / \sqrt{\pi}$, we find the charge $\Delta\theta/2\pi$, as from our earlier derivation.

Although the σ model proper does not support finite-energy solitons, we can consider a fermion interacting with external fields of this type. This proves useful as a warmup for the gauge theory monopoles to be discussed shortly.

The interaction Lagrangian is of standard form

$$\mathcal{L}_I = g \bar{\psi} (\varphi_0 + i \bar{\psi} \cdot \vec{\tau} \gamma_5) \psi$$

with ψ an isodoublet fermion field. We compute the induced current as in the 1+1 dimensional examples, from graphs as in Fig. 3. A straight-

forward calculation leads to

$$\langle j^\mu \rangle = \frac{1}{12\pi^2 |\psi|^4} \epsilon^{\mu\alpha\beta\gamma} \epsilon_{abcd} \varphi_a \partial_\alpha \varphi_b \partial_\beta \varphi_c \partial_\gamma \varphi_d. \quad (6)$$

With this form, $\partial_\mu \langle j^\mu \rangle = 0$. This, of course, indicates that only the behavior at spatial infinity determines the charge, since changes in the fields in a finite volume lead only to current flows in a finite volume and therefore do not change the total charge. In fact, if we take $\varphi_0 = m/g$, $\varphi_a = \hat{\varphi}_a(\vec{x}) f(t)$, $a = 1, 2, 3$, where $\hat{\varphi}_a(\vec{x}) = vx_a/|x|$ as $|x| \rightarrow \infty$; and evaluate the current

$$\langle j^\mu \rangle = \frac{1}{12\pi^2} \epsilon^{\mu\alpha\beta\gamma} \epsilon_{abcd} \left[\frac{\varphi_0}{|\psi|^4} (\nabla_\alpha \varphi)_a (\nabla_\beta \varphi)_b (\nabla_\gamma \varphi)_c + \frac{3}{4} e F_{\alpha\beta,ab} \frac{\varphi_0}{|\psi|^2} (\nabla_\gamma \varphi)_c \right] \quad (8)$$

obeys $\partial_\mu \langle j^\mu \rangle = (-e^2/128\pi^2) \epsilon^{\alpha\beta\gamma\delta} \epsilon_{abcd} F_{\alpha\beta,ab} F_{\gamma\delta,cd}$. This is the expected anomaly and vanishes when we have only vector gauge fields as in the monopole. The coefficient of the second term in (8) can be checked by the evaluation of the diagram in Fig. 3 with one gauge field vertex inserted.

We now take φ as before and $A_{ab} = \hat{A}_{ab}(\vec{x})$, $A_{a0} = 0$, $a, b = 1, 2, 3$, where $\hat{\varphi}$ and \hat{A} are the monopole fields, and find the current flow at infinity. Since $(\nabla_\mu \varphi)_a = 0$ at infinity, the only contribution comes from taking $\gamma = d = 0$ in the second term of (8) and gives for the fermion number

$$(e\Phi/4\pi^2) \tan^{-1}(gv/m), \quad (9)$$

where Φ is the magnetic flux out of the sphere at infinity. Since $e\Phi = 4\pi$, this gives fermion number $\frac{1}{2}$ when $m \rightarrow 0$!

The direct utility of our results for particle physics is highly problematical. Even if magnetic monopoles were found, their fermion number is not a reasonable quantity in standard theories. [In principle, we could imagine coupling a U(1) gauge field to the fermion number, and so the calculation is not entirely content free!] We do think that the results are an interesting curiosity in quantum field theory and as such may eventually be useful. It is likely that kindred, but experimentally accessible, effects do arise in condensed matter systems.

We are especially grateful to J. R. Schrieffer for interesting us in this problem and to L. Susskind for reminding us of the bosonization method.

flow at infinity, we find a fermion number

$$\pi^{-1}(\theta - \sin\theta \cos\theta), \quad \tan\theta = gv/m, \quad (7)$$

which $\rightarrow \frac{1}{2}$ as $m \rightarrow 0$.

We may extend this analysis in a simple way to the monopole solutions of non-Abelian gauge theories by simply gauging the $SU(2) \otimes SU(2)$ chiral symmetry of our σ model. In the end, we can specialize by setting the axial gauge fields to zero, and fixing a fermion mass ($\varphi_0 = \text{const}$).

The expression (6) for the current is changed in the first instance by the conversion of ordinary to covariant derivatives, $\partial \rightarrow \nabla \equiv \partial + eA$. This is not sufficient, however, since this minimally modified current is not conserved. The current

This work was started while the first author was at the Santa Barbara Institute. This work was supported in part by the U. S. Department of Energy under Contract No. DE-AC03-SF00515 and by the National Science Foundation under Grant No. PHY 77-27084.

^(a)Permanent address: Center for Theoretical Physics, Massachusetts Institute of Technology, Cambridge, Mass. 02139

¹R. Jackiw and C. Rebbi, Phys. Rev. D 13, 3398 (1978).

²W. P. Su, J. R. Schrieffer, and A. J. Heeger, Phys. Rev. Lett. 42, 1698 (1979), and Phys. Rev. B 22, 2099 (1980).

³R. Jackiw and J. R. Schrieffer, Santa Barbara Report No. NSF-ITP-81-01 (to be published).

⁴P. Hasenfratz and G. 't Hooft, Phys. Rev. Lett. 36, 1119 (1976).

⁵R. Jackiw and C. Rebbi, Phys. Rev. Lett. 36, 1116 (1976).

⁶F. Wilczek, to be published.

⁷E. Witten, Phys. Lett. 86B, 283 (1979).

⁸F. Wilczek, to be published.

⁹For a discussion of both theoretical and experimental aspects of polyacetylene, see A. Heeger, Comments Solid State Phys. 10, 53 (1981).

¹⁰W. P. Su and J. R. Schrieffer, Phys. Rev. Lett. 46, 738 (1981).

¹¹See, e.g., S. Mandelstam, Phys. Rev. D 11, 3026 (1975); S. Coleman, R. Jackiw, and L. Susskind, Ann. Phys. (N.Y.) 93, 267 (1975).

GLOBAL ASPECTS OF CURRENT ALGEBRA

Edward WITTEN*

Joseph Henry Laboratories, Princeton University, Princeton, New Jersey 08544, USA

Received 4 March 1983

A new mathematical framework for the Wess-Zumino chiral effective action is described. It is shown that this action obeys an a priori quantization law, analogous to Dirac's quantization of magnetic charge. It incorporates in current algebra both perturbative and non-perturbative anomalies.

The purpose of this paper is to clarify an old but relatively obscure aspect of current algebra: the Wess-Zumino effective lagrangian [1] which summarizes the effects of anomalies in current algebra. As we will see, this effective lagrangian has unexpected analogies to some 2 + 1 dimensional models discussed recently by Deser et al. [2] and to a recently noted SU(2) anomaly [3]. There also are connections with work of Balachandran et al. [4].

For definiteness we will consider a theory with $SU(3)_L \times SU(3)_R$ symmetry spontaneously broken down to the diagonal $SU(3)$. We will ignore explicit symmetry-breaking perturbations, such as quark bare masses. With $SU(3)_L \times SU(3)_R$ broken to diagonal $SU(3)$, the vacuum states of the theory are in one to one correspondence with points in the $SU(3)$ manifold. Correspondingly, the low-energy dynamics can be conveniently described by introducing a field $U(x^\alpha)$ that transforms in a so-called non-linear realization of $SU(3)_L \times SU(3)_R$. For each space-time point x^α , $U(x^\alpha)$ is an element of $SU(3)$: a 3×3 unitary matrix of determinant one. Under an $SU(3)_L \times SU(3)_R$ transformation by unitary matrices (A, B) , U transforms as $U \rightarrow AUB^{-1}$.

The effective lagrangian for U must have $SU(3)_L \times SU(3)_R$ symmetry, and, to describe correctly the low-energy limit, it must have the smallest possible number of derivatives. The unique choice with only two derivatives is

$$\mathcal{L} = \frac{1}{16} F_\pi^2 \int d^4x \text{Tr } \partial_\mu U \partial_\mu U^{-1}, \quad (1)$$

* Supported in part by NSF Grant PHY80-19754.

where experiment indicates $F_\pi \approx 190$ MeV. The perturbative expansion of U is

$$U = 1 + \frac{2i}{F_\pi} \sum_{a=1}^8 \lambda^a \pi^a + \dots, \quad (2)$$

where λ^a (normalized so $\text{Tr } \lambda^a \lambda^b = 2\delta^{ab}$) are the SU(3) generators and π^a are the Goldstone boson fields.

This effective lagrangian is known to incorporate all relevant symmetries of QCD. All current algebra theorems governing the extreme low-energy limit of Goldstone boson S -matrix elements can be recovered from the tree approximation to it. What is less well known, perhaps, is that (1) possesses an extra discrete symmetry that is *not* a symmetry of QCD.

The lagrangian (1) is invariant under $U \leftrightarrow U^T$. In terms of pions this is $\pi^0 \leftrightarrow \pi^0$, $\pi^+ \leftrightarrow \pi^-$; it is ordinary charge conjugation. (1) is also invariant under the naive parity operation $x \leftrightarrow -x$, $t \leftrightarrow t$, $U \leftrightarrow U$. We will call this P_0 . And finally, (1) is invariant under $U \leftrightarrow U^{-1}$. Comparing with eq. (2), we see that this latter operation is equivalent to $\pi^a \leftrightarrow -\pi^a$, $a = 1, \dots, 8$. This is the operation that counts modulo two the number of bosons, N_B , so we will call it $(-1)^{N_B}$.

Certainly, $(-1)^{N_B}$ is not a symmetry of QCD. The problem is the following. QCD is parity invariant only if the Goldstone bosons are treated as pseudoscalars. The parity operation in QCD corresponds to $x \leftrightarrow -x$, $t \leftrightarrow t$, $U \leftrightarrow U^{-1}$. This is $P = P_0(-1)^{N_B}$. QCD is invariant under P but not under P_0 or $(-1)^{N_B}$ separately. The simplest process that respects all bona fide symmetries of QCD but violates P_0 and $(-1)^{N_B}$ is $K^+ K^- \rightarrow \pi^+ \pi^0 \pi^-$ (note that the ϕ meson decays to both $K^+ K^-$ and $\pi^+ \pi^0 \pi^-$). It is natural to ask whether there is a simple way to add a higher-order term to (1) to obtain a lagrangian that obeys *only* the appropriate symmetries.

The Euler-Lagrangian equation derived from (1) can be written

$$\partial_\mu \left(\frac{1}{8} F_\pi^2 U^{-1} \partial_\mu U \right) = 0. \quad (3)$$

Let us try to add a suitable extra term to this equation. A Lorentz-invariant term that violates P_0 must contain the Levi-Civita symbol $\epsilon_{\mu\nu\alpha\beta}$. In the spirit of current algebra, we wish a term with the smallest possible number of derivatives, since, in the low-energy limit, the derivatives of U are small. There is a unique P_0 -violating term with only four derivatives. We can generalize (3) to

$$\partial_\mu \left(\frac{1}{8} F_\pi^2 U^{-1} \partial_\mu U \right) + \lambda \epsilon^{\mu\nu\alpha\beta} U^{-1} (\partial_\mu U) U^{-1} (\partial_\nu U) U^{-1} (\partial_\alpha U) U^{-1} (\partial_\beta U) = 0, \quad (4)$$

λ being a constant. Although it violates P_0 , (4) can be seen to respect $P = P_0(-1)^{N_B}$.

Can eq. (4) be derived from a lagrangian? Here we find trouble. The only pseudoscalar of dimension four would seem to be $\epsilon^{\mu\nu\alpha\beta} \text{Tr } U^{-1} (\partial_\mu U) \cdot U^{-1} (\partial_\nu U) U^{-1} (\partial_\alpha U) U^{-1} (\partial_\beta U)$, but this vanishes, by antisymmetry of $\epsilon^{\mu\nu\alpha\beta}$ and cyclic symmetry of the trace. Nevertheless, as we will see, there is a lagrangian.

Let us consider a simple problem of the same sort. Consider a particle of mass m constrained to move on an ordinary two-dimensional sphere of radius one. The lagrangian is $\mathcal{L} = \frac{1}{2}m\int dt \dot{x}_i^2$ and the equation of motion is $m\ddot{x}_i + mx_i(\sum_k \dot{x}_k^2) = 0$; the constraint is $\sum x_i^2 = 1$. This system respects the symmetries $t \leftrightarrow -t$ and separately $x_i \leftrightarrow -x_i$. If we want an equation that is only invariant under the combined operation $t \leftrightarrow -t$, $x_i \leftrightarrow x_i$, the simplest choice is

$$m\ddot{x}_i + mx_i \left(\sum_k \dot{x}_k^2 \right) = \alpha \epsilon_{ijk} x_j \dot{x}_k, \quad (5)$$

where α is a constant. To derive this equation from a lagrangian is again troublesome. There is no obvious term whose variation equals the right-hand side (since $\epsilon_{ijk} x_i x_j \dot{x}_k = 0$).

However, this problem has a well-known solution. The right-hand side of (5) can be understood as the Lorentz force for an electric charge interacting with a magnetic monopole located at the center of the sphere. Introducing a vector potential A such that $\nabla \times A = x/|x|^3$, the action for our problem is

$$I = \int \left(\frac{1}{2}m\dot{x}_i^2 + \alpha A_i \dot{x}_i \right) dt. \quad (6)$$

This lagrangian is problematical because A_i contains a Dirac string and certainly does not respect the symmetries of our problem. To explore this quantum mechanically let us consider the simplest form of the Feynman path integral, $\text{Tr } e^{-\beta H} = \int dx_i(t) e^{-S}$. In e^{-S} the troublesome term is

$$\exp \left(i\alpha \int_{\gamma} A_i dx' \right), \quad (7)$$

where the integration goes over the particle orbit γ : a closed orbit if we discuss the simplest object $\text{Tr } e^{-\beta H}$.

By Gauss's law we can eliminate the vector potential from (7) in favor of the magnetic field. In fact, the closed orbit γ of fig. 1a is the boundary of a disc D , and by Gauss's law we can write (7) in terms of the magnetic flux through D :

$$\exp \left(i\alpha \int_{\gamma} A_i dx' \right) = \exp \left(i\alpha \int_D F_{ij} d\Sigma' \right). \quad (8)$$

The precise mathematical statement here is that since $\pi_1(S^2) = 0$, the circle γ in S^2 is the boundary of a disc D (or more exactly, a mapping γ of a circle into S^2 can be extended to a mapping of a disc into S^2).

The right-hand side of (8) is manifestly well defined, unlike the left-hand side, which suffers from a Dirac string. We could try to use the right-hand side of (8) in a Feynman path integral. There is only one problem: D isn't unique. The curve γ also bounds the disc D' (fig. 1c). There is no consistent way to decide whether to choose

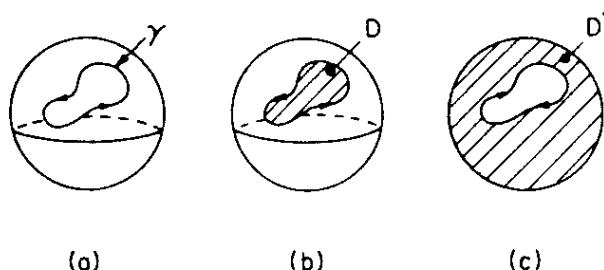


Fig. 1. A particle orbit γ on the two-sphere (part (a)) bounds the discs D (part (b)) and D' (part (c)).

D or D' (the curve γ could continuously be looped around the sphere or turned inside out). Working with D' we would get

$$\exp\left(i\alpha \int_{\gamma} A_i dx^i\right) = \exp\left(-i\alpha \int_{D'} F_{ij} d\Sigma^{ij}\right), \quad (9)$$

where a crucial minus sign on the right-hand side of (9) appears because γ bounds D in a right-hand sense, but bounds D' in a left-hand sense. If we are to introduce the right-hand side of (8) or (9) in a Feynman path integral, we must require that they be equal. This is equivalent to

$$1 = \exp\left(i\alpha \int_{D+D'} F_{ij} d\Sigma^{ij}\right). \quad (10)$$

Since $D + D'$ is the whole two sphere S^2 , and $\int_{S^2} F_{ij} d\Sigma^{ij} = 4\pi$, (10) is obeyed if and only if α is an integer or half-integer. This is Dirac's quantization condition for the product of electric and magnetic charges.

Now let us return to our original problem. We imagine space-time to be a very large four-dimensional sphere M . A given non-linear sigma model field U is a mapping of M into the $SU(3)$ manifold (fig. 2a). Since $\pi_4(SU(3)) = 0$, the four-sphere in $SU(3)$ defined by $U(x)$ is the boundary of a five-dimensional disc Q .

By analogy with the previous problem, let us try to find some object that can be integrated over Q to define an action functional. On the $SU(3)$ manifold there is a unique fifth rank antisymmetric tensor ω_{ijklm} that is invariant under $SU(3)_L \times SU(3)_R$ *. Analogous to the right-hand side of eq. (8), we define

$$\Gamma = \int_Q \omega_{ijklm} d\Sigma^{ijklm}. \quad (11)$$

* Let us first try to define ω at $U=1$; it can then be extended to the whole $SU(3)$ manifold by an $SU(3)_L \times SU(3)_R$ transformation. At $U=1$, ω must be invariant under the diagonal subgroup of $SU(3)_L \times SU(3)_R$ that leaves fixed $U=1$. The tangent space to the $SU(3)$ manifold at $U=1$ can be identified with the Lie algebra of $SU(3)$. So ω , at $U=1$, defines a fifth-order antisymmetric invariant in the $SU(3)$ Lie algebra. There is only one such invariant. Given five $SU(3)$ generators A, B, C, D and E , the one such invariant is $\text{Tr } ABCDE - \text{Tr } BACDE \pm \text{permutations}$. The $SU(3)_L \times SU(3)_R$ invariant ω so defined has zero curl ($\partial_{[i} \omega_{jklm]} \pm \text{permutations} = 0$) and for this reason (11) is invariant under infinitesimal variations of Q ; there arises only the topological problem discussed in the text.

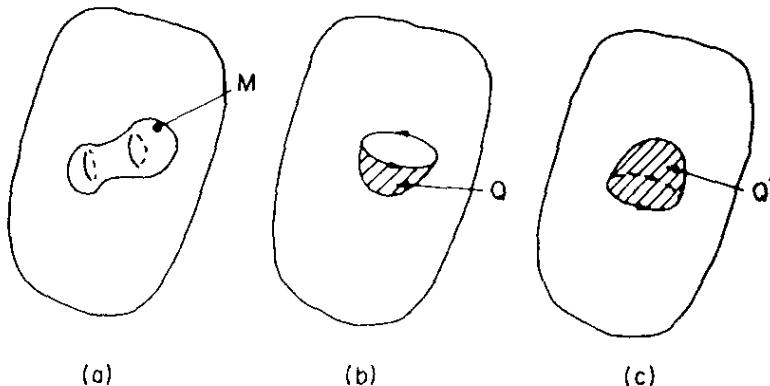


Fig. 2. Space-time, a four-sphere, is mapped into the $SU(3)$ manifold. In part (a), space-time is symbolically denoted as a two-sphere. In parts (b) and (c), space-time is reduced to a circle that bounds the discs Q and Q' . The $SU(3)$ manifold is symbolized in these sketches by the interior of the oblong.

As before, we hope to include $\exp(i\Gamma)$ in a Feynman path integral. Again, the problem is that Q is not unique. Our four-sphere M is also the boundary of another five-disc Q' (fig. 2c). If we let

$$\Gamma' = - \int_{Q'} \omega_{ijklm} d\Sigma^{ijklm}, \quad (12)$$

(with, again, a minus sign because M bounds Q' with opposite orientation) then we must require $\exp(i\Gamma) = \exp(i\Gamma')$ or equivalently $\int_{Q+Q'} \omega_{ijklm} d\Sigma^{ijklm} = 2\pi \cdot \text{integer}$. Since $Q + Q'$ is a closed five-dimensional sphere, our requirement is

$$\int_S \omega_{ijklm} d\Sigma^{ijklm} = 2\pi \cdot \text{integer},$$

for any five-sphere S in the $SU(3)$ manifold.

We thus need the topological classification of mappings of the five-sphere into $SU(3)$. Since $\pi_5(SU(3)) = \mathbb{Z}$, every five sphere in $SU(3)$ is topologically a multiple of a basic five sphere S_0 . We normalize ω so that

$$\int_{S_0} \omega_{ijklm} d\Sigma^{ijklm} = 2\pi, \quad (13)$$

and then (with Γ in eq. (11)) we may work with the action

$$I = \frac{1}{16} F_\pi^2 \int d^4x \operatorname{Tr} \partial_\mu U \partial_\mu U^{-1} + n\Gamma, \quad (14)$$

where n is an arbitrary integer. Γ is, in fact, the Wess-Zumino lagrangian. Only the a priori quantization of n is a new result.

The identification of S_0 and the proper normalization of ω is a subtle mathematical problem. The solution involves a factor of two from the Bott periodicity theorem. Without abstract notation, the result [5] can be stated as follows. Let $y^i, i = 1 \dots 5$ be coordinates for the disc Q . Then on Q (where we need it)

$$d\Sigma^{ijklm} \omega_{ijklm} = -\frac{i}{240\pi^2} d\Sigma^{ijklm} \left[\text{Tr } U^{-1} \frac{\partial U}{\partial y^i} U^{-1} \frac{\partial U}{\partial y^j} U^{-1} \frac{\partial U}{\partial y^k} U^{-1} \frac{\partial U}{\partial y^l} U^{-1} \frac{\partial U}{\partial y^m} \right]. \quad (15)$$

The physical consequences of this can be made more transparent as follows. From eq. (2),

$$U^{-1} \partial_i U = \frac{2i}{F_\pi} \partial_i A + O(A^2), \quad \text{where } A = \Sigma \lambda^a \pi^a. \quad (16)$$

So

$$\begin{aligned} \omega_{ijklm} d\Sigma^{ijklm} &= \frac{2}{15\pi^2 F_\pi^5} d\Sigma^{ijklm} \text{Tr } \partial_i A \partial_j A \partial_k A \partial_l A \partial_m A + O(A^6) \\ &= \frac{2}{15\pi^2 F_\pi^5} d\Sigma^{ijklm} \partial_i (\text{Tr } A \partial_j A \partial_k A \partial_l A \partial_m A) + O(A^6). \end{aligned}$$

So $\int_Q \omega_{ijklm} d\Sigma^{ijklm}$ is (to order A^5 and in fact also in higher orders) the integral of a total divergence which can be expressed by Stokes' theorem as an integral over the boundary of Q . By construction, this boundary is precisely space-time. We have, then,

$$n\Gamma = n \frac{2}{15\pi^2 F_\pi^5} \int d^4x \epsilon^{\mu\nu\alpha\beta} \text{Tr } A \partial_\mu A \partial_\nu A \partial_\alpha A \partial_\beta A + \text{higher order terms}. \quad (17)$$

In a hypothetical world of massless kaons and pions, this effective lagrangian rigorously describes the low-energy limit of $K^+ K^- \rightarrow \pi^+ \pi^0 \pi^-$ ^{*}. We reach the remarkable conclusion that in any theory with $SU(3) \times SU(3)$ broken to diagonal $SU(3)$, the low-energy limit of the amplitude for this reaction must be (in units given in (17)) an integer.

What is the value of this integer in QCD? Were n to vanish, the practical interest of our discussion would be greatly reduced. It turns out that if N_c is the number of colors (three in the real world) then $n = N_c$. The simplest way to deduce this is a

* Our formula should agree for $n = 1$ with formulas of ref. [1], as later equations make clear. There appears to be a numerical error on p. 97 of ref. [1] ($\frac{1}{6}$ instead of $\frac{2}{15}$).

procedure that is of interest anyway, viz. coupling to electromagnetism, so as to describe the low-energy dynamics of Goldstone bosons and photons.

Let

$$Q = \begin{pmatrix} \frac{2}{3} & & \\ & -\frac{1}{3} & \\ & & -\frac{1}{3} \end{pmatrix}$$

be the usual electric charge matrix of quarks. The functional Γ is invariant under global charge rotations, $U \rightarrow U + i\varepsilon[Q, U]$, where ε is a constant. We wish to promote this to a local symmetry, $U \rightarrow U + i\varepsilon(x)[Q, U]$, where $\varepsilon(x)$ is an arbitrary function of x . It is necessary, of course, to introduce the photon field A_μ which transforms as $A_\mu \rightarrow A_\mu - (1/e)\partial_\mu\varepsilon$; e is the charge of the proton.

Usually a global symmetry can straightforwardly be gauged by replacing derivatives by covariant derivatives, $\partial_\mu \rightarrow D_\mu = \partial_\mu + ieA_\mu$. In the case at hand, Γ is not given as the integral of a manifestly $SU(3)_L \times SU(3)_R$ invariant expression, so the standard road to gauging global symmetries of Γ is not available. One can still resort to the trial and error Noether method, widely used in supergravity. Under a local charge rotation, one finds $\Gamma \rightarrow \Gamma - \int d^4x \partial_\mu\varepsilon J^\mu$ where

$$\begin{aligned} J^\mu = \frac{1}{48\pi^2} \epsilon^{\mu\nu\alpha\beta} \text{Tr} [& Q(\partial_\nu U U^{-1})(\partial_\alpha U U^{-1})(\partial_\beta U U^{-1}) \\ & + Q(U^{-1}\partial_\nu U)(U^{-1}\partial_\alpha U)(U^{-1}\partial_\beta U)] , \end{aligned} \quad (18)$$

is the extra term in the electromagnetic current required (from Noether's theorem) due to the addition of Γ to the lagrangian. The first step in the construction of an invariant lagrangian is to add the Noether coupling, $\Gamma \rightarrow \Gamma' = \Gamma - e \int d^4x A_\mu J^\mu(x)$. This expression is still not gauge invariant, because J^μ is not, but by trial and error one finds that by adding an extra term one can form a gauge invariant functional

$$\begin{aligned} \tilde{\Gamma}(U, A_\mu) = \Gamma(U) - e \int d^4x A_\mu J^\mu + \frac{ie^2}{24\pi^2} \int d^4x \epsilon^{\mu\nu\alpha\beta} (\partial_\mu A_\nu) A_\alpha \\ \times \text{Tr} [Q^2(\partial_\beta U) U^{-1} + Q^2 U^{-1}(\partial_\beta U) + Q U Q U^{-1}(\partial_\beta U) U^{-1}] . \end{aligned} \quad (19)$$

Our gauge invariant lagrangian will then be

$$\mathcal{L} = \frac{1}{16} F_\pi^2 \int d^4x \text{Tr} D_\mu U D_\mu U^{-1} + n \tilde{\Gamma} . \quad (20)$$

What value of the integer n will reproduce QCD results?

Here we find a surprise. The last term in (18) has a piece that describes $\pi^0 \rightarrow \gamma\gamma$. Expanding U and integrating by parts, (18) has a piece

$$A = \frac{ne^2}{48\pi^2 F_\pi} \pi^0 \epsilon^{\mu\nu\alpha\beta} F_{\mu\nu} F_{\alpha\beta}. \quad (21)$$

This agrees with the result from QCD triangle diagrams [6] if $n = N_c$, the number of colors. The Noether coupling $-e A_\mu J^\mu$ describes, among other things, a $\gamma\pi^+ \pi^0 \pi^-$ vertex

$$B = -\frac{2}{3}ie \frac{n}{\pi^2 F_\pi^3} \epsilon^{\mu\nu\alpha\beta} A_\mu \partial_\nu \pi^+ \partial_\alpha \pi^- \partial_\beta \pi^0. \quad (22)$$

Again this agrees with calculations [7] based on the QCD VAAA anomaly if $n = N_c$. The effective action $N_c \tilde{\Gamma}$ (first constructed in another way by Wess and Zumino) precisely describes all effects of QCD anomalies in low-energy processes with photons and Goldstone bosons.

It is interesting to try to gauge subgroups of $SU(3)_L \times SU(3)_R$ other than electromagnetism. One may have in mind, for instance, applications to the standard weak interaction model. In general, one may try to gauge an arbitrary subgroup H of $SU(3)_L \times SU(3)_R$, with generators K^σ , $\sigma = 1 \dots r$. Each K^σ is a linear combination of generators T_L^σ and T_R^σ of $SU(3)_L$ and $SU(3)_R$, $K^\sigma = T_L^\sigma + T_R^\sigma$. (Either T_L^σ or T_R^σ may vanish for some values of σ .) For any space-time dependent functions $\epsilon^\sigma(x)$, let $\epsilon_L = \sum_\sigma T_L^\sigma \epsilon^\sigma(x)$, $\epsilon_R = \sum_\sigma T_R^\sigma \epsilon^\sigma(x)$. We want an action with local invariance under $U \rightarrow U + i(\epsilon_L(x)U - U\epsilon_R(x))$.

Naturally, it is necessary to introduce gauge fields $A_\mu^\sigma(x)$, transforming as $A_\mu^\sigma(x) \rightarrow A_\mu^\sigma(x) - (1/e_\sigma) \partial_\mu \epsilon^\sigma + f^{\sigma\rho\sigma} \epsilon^\rho A_\mu^\rho$ where e_σ is the coupling constant corresponding to the generator K^σ , and $f^{\sigma\rho\sigma}$ are the structure constants of H . It is useful to define $A_{\mu L} = \sum_\sigma e_\sigma A_\mu^\sigma T_L^\sigma$, $A_{\mu R}^R = \sum_\sigma e_\sigma A_\mu^\sigma T_R^\sigma$.

We have already seen that Γ incorporates the effects of anomalies, so it is not very surprising that a generalization of Γ that is gauge invariant under H exists only if H is a so-called anomaly-free subgroup of $SU(3)_L \times SU(3)_R$. Specifically, one finds that H can be gauged only if for each σ ,

$$\text{Tr}(T_L^\sigma)^3 = \text{Tr}(T_R^\sigma)^3, \quad (23)$$

which is the usual condition for cancellation of anomalies at the quark level.

If (23) is obeyed, a gauge invariant generalization of Γ can be constructed somewhat tediously by trial and error. It is useful to define $U_{\mu L} = (\partial_\nu U) U^{-1}$ and $U_{\nu R} = U^{-1} \partial_\nu U$. The gauge invariant functional then turns out to be

$$\tilde{\Gamma}(A_\mu, U) = \Gamma(U) + \frac{1}{48\pi^2} \int d^4x \epsilon^{\mu\nu\alpha\beta} Z_{\mu\nu\alpha\beta},$$

where

$$\begin{aligned}
Z_{\mu\nu\alpha\beta} = & - \text{Tr} [A_{\mu L} U_{\nu L} U_{\alpha L} U_{\beta L} + (L \rightarrow R)] \\
& + i \text{Tr} [[(\partial_\mu A_{\nu L}) A_{\alpha L} + A_{\mu L} (\partial_\nu A_{\alpha L})] U_{\beta L} + (L \rightarrow R)] \\
& + i \text{Tr} [(\partial_\mu A_{\nu R}) U^{-1} A_{\alpha L} \partial_\beta U + A_{\mu L} U^{-1} (\partial_\nu A_{\alpha R}) \partial_\beta U] \\
& - \frac{1}{2} i \text{Tr} (A_{\mu L} U_{\nu L} A_{\alpha L} U_{\beta L} - (L \rightarrow R)) \\
& + i \text{Tr} [A_{\mu L} U A_{\nu R} U^{-1} U_{\alpha L} U_{\beta L} - A_{\mu R} U^{-1} A_{\nu L} U U_{\alpha R} U_{\beta R}] \\
& - \text{Tr} [[(\partial_\mu A_{\nu R}) A_{\alpha R} + A_{\mu R} (\partial_\nu A_{\alpha R})] U^{-1} A_{\beta L} U \\
& - [(\partial_\mu A_{\nu L}) A_{\alpha L} + A_{\mu L} (\partial_\nu A_{\alpha L})] U A_{\beta R} U^{-1}] \\
& - \text{Tr} [A_{\mu R} U^{-1} A_{\nu L} U A_{\alpha R} U_{\beta R} + A_{\mu L} U A_{\nu R} U^{-1} A_{\alpha L} U_{\beta L}] \\
& - \text{Tr} [A_{\mu L} A_{\nu L} U (\partial_\alpha A_{\beta R}) U^{-1} + A_{\mu R} A_{\nu R} U^{-1} (\partial_\alpha A_{\beta L}) U] \\
& - i \text{Tr} [A_{\mu R} A_{\nu R} A_{\alpha R} U^{-1} A_{\beta L} U - A_{\mu L} A_{\nu L} A_{\alpha L} U A_{\beta R} U^{-1} \\
& + \frac{1}{2} A_{\mu L} A_{\nu L} U A_{\alpha R} A_{\beta R} U^{-1} + \frac{1}{2} A_{\mu R} U^{-1} A_{\nu L} U A_{\alpha R} U^{-1} A_{\beta L} U]. \quad (24)
\end{aligned}$$

If eq. (22) for cancellation of anomalies is not obeyed, then the variation of $\tilde{\Gamma}$ under a gauge transformation does not vanish but is

$$\begin{aligned}
\delta \tilde{\Gamma} = & - \frac{1}{24\pi^2} \int d^4x \epsilon^{\mu\nu\alpha\beta} \text{Tr} \epsilon_L [(\partial_\mu A_{\nu L}) (\partial_\alpha A_{\beta L}) - \frac{1}{2} i \partial_\mu (A_{\nu L} A_{\alpha L} A_{\beta L})] \\
& - (L \rightarrow R), \quad (25)
\end{aligned}$$

in agreement with computations at the quark level [8] of the anomalous variation of the effective action under a gauge transformation.

Thus, Γ incorporates all information usually associated with triangle anomalies, including the restriction on what subgroups H of $SU(3)_L \times SU(3)_R$ can be gauged. However, there is another potential obstruction to the ability to gauge a subgroup of $SU(3)_L \times SU(3)_R$. This is the non-perturbative anomaly [3] associated with $\pi_4(H)$. Is this anomaly, as well, implicit in Γ ? In fact, it is.

Let H be an $SU(2)$ subgroup of $SU(3)_L$, chosen so that an $SU(2)$ matrix W is embedded in $SU(3)_L$ as

$$\hat{W} = \left(\begin{array}{c|c} W & 0 \\ \hline 0 & 0 \end{array} \right).$$

This subgroup is free of triangle anomalies, so the functional $\tilde{\Gamma}$ of eq. (23) is invariant under infinitesimal local H transformations.

However, is $\tilde{\Gamma}$ invariant under H transformations that cannot be reached continuously? Since $\pi_4(\text{SU}(2)) = \mathbb{Z}_2$, there is one non-trivial homotopy class of SU(2) gauge transformations. Let W be an SU(2) gauge transformation in this non-trivial class. Under \hat{W} , $\tilde{\Gamma}$ may at most be shifted by a constant, independent of U and A_μ , because $\delta\tilde{\Gamma}/\delta U$ and $\delta\tilde{\Gamma}/\delta A_\mu$ are gauge-covariant local functionals of U and A_μ . Also $\tilde{\Gamma}$ is invariant under \hat{W}^2 , since \hat{W}^2 is equivalent to the identity in $\pi_4(\text{SU}(2))$, and we know $\tilde{\Gamma}$ is invariant under topologically trivial gauge transformations. This does not quite mean that $\tilde{\Gamma}$ is invariant under W . Since $\tilde{\Gamma}$ is only defined modulo 2π , the fact that $\tilde{\Gamma}$ is invariant under W^2 leaves two possibilities for how $\tilde{\Gamma}$ behaves under W . It may be invariant, or it may be shifted by π .

To choose between these alternatives, it is enough to consider a special case. For instance, it suffices to evaluate $\Delta = \tilde{\Gamma}(U = 1, A_\mu = 0) - \tilde{\Gamma}(U = \hat{W}, A_\mu = ie^{-1}(\partial_\mu \hat{W})\hat{W}^{-1})$. It is not difficult to see that in this case the complicated terms involving $\epsilon^{\mu\nu\alpha\beta} Z_{\mu\nu\alpha\beta}$ vanish, so in fact $\Delta = \Gamma(U = 1) - \Gamma(U = \hat{W})$. A detailed calculation shows that

$$\Gamma(U = 1) - \Gamma(U = \hat{W}) = \pi. \quad (26)$$

This calculation has some other interesting applications and will be described elsewhere [9].

The Feynman path integral, which contains a factor $\exp(iN_c \tilde{\Gamma})$, hence picks up under W a factor $\exp(iN_c \pi) = (-1)^{N_c}$. It is gauge invariant if N_c is even, but not if N_c is odd. This agrees with the determination of the SU(2) anomaly at the quark level [3]. For under H, the right-handed quarks are singlets. The left-handed quarks consist of one singlet and one doublet per color, so the number of doublets equals N_c . The argument of ref. [3] shows at the quark level that the effective action transforms under W as $(-1)^{N_c}$.

Finally, let us make the following remark, which apart from its intrinsic interest will be useful elsewhere [9]. Consider $\text{SU}(3)_L \times \text{SU}(3)_R$ currents defined at the quark level as

$$J_{\mu L}^a = \bar{q} \lambda^a \gamma_{\mu}^{\frac{1}{2}} (1 - \gamma_5) q, \quad J_{\mu R}^a = \bar{q} \lambda^a \gamma_{\mu}^{\frac{1}{2}} (1 + \gamma_5) q. \quad (27)$$

By analogy with eq. (17), the proper sigma model description of these currents contains pieces

$$J_L^{\mu a} = \frac{N_c}{48\pi^2} \epsilon^{\mu\nu\alpha\beta} \text{Tr} \lambda^a U_{\nu L} U_{\alpha L} U_{\beta L},$$

$$J_R^{\mu a} = \frac{N_c}{48\pi^2} \epsilon^{\mu\nu\alpha\beta} \text{Tr} \lambda^a U_{\nu R} U_{\alpha R} U_{\beta R}, \quad (28)$$

corresponding (via Noether's theorem) to the addition to the lagrangian of $N_c \Gamma$. In this discussion, the λ^α should be traceless SU(3) generators. However, let us try to construct an anomalous baryon number current in the same way. We define the baryon number of a quark (whether left-handed or right-handed) to be $1/N_c$, so that an ordinary baryon made from N_c quarks has baryon number one. Replacing λ^α by $1/N_c$, but including contributions of both left-handed and right-handed quarks, the anomalous baryon-number current would be

$$J^\mu = \frac{1}{24\pi^2} \epsilon^{\mu\nu\alpha\beta} \text{Tr } U^{-1} \partial_\nu U U^{-1} \partial_\alpha U U^{-1} \partial_\beta U. \quad (29)$$

One way to see that this is the proper, and properly normalized, formula is to consider gauging an arbitrary subgroup not of $SU(3)_L \times SU(3)_R$ but of $SU(3)_L \times SU(3)_R \times U(1)$, $U(1)$ being baryon number. The gauging of $U(1)$ is accomplished by adding a Noether coupling $-e J^\mu B_\mu$ plus whatever higher-order terms may be required by gauge invariance. (B_μ is a $U(1)$ gauge field which may be coupled as well to some $SU(3)_L \times SU(3)_R$ generator.) With J^μ defined in (29), this leads to a generalization of $\tilde{\Gamma}$ that properly reflects anomalous diagrams involving the baryon-number current (for instance, it properly incorporates the anomaly in the baryon number $SU(2)_L - SU(2)_L$ triangle that leads to baryon non-conservation by instantons in the standard weak interaction model). Eq. (29) may also be extracted from QCD by methods of Goldstone and Wilczek [10].

References

- [1] J. Wess and B. Zumino, Phys. Lett. 37B (1971) 95
- [2] S. Deser, R. Jackiw and S. Templeton, Phys. Rev. Lett. 48 (1982) 975; Ann. of Phys. 140 (1982) 372
- [3] E. Witten, Phys. Lett. 117B (1982) 324
- [4] A.P. Balachandran, V.P. Nair and C.G. Trahern, Syracuse University preprint SU-4217-205 (1981)
- [5] R. Bolt and R. Seeley, Comm. Math. Phys. 62 (1978) 235
- [6] S.L. Adler, Phys. Rev. 177 (1969) 2426;
J.S. Bell and R. Jackiw, Nuovo Cim. 60 (1969) 147;
W.A. Bardeen, Phys. Rev. 184 (1969) 1848
- [7] S.L. Adler and W.A. Bardeen, Phys. Rev. 182 (1969) 1517;
R. Aviv and A. Zee, Phys. Rev. D5 (1972) 2372
S.L. Adler, B.W. Lee, S.B. Treiman and A. Zee, Phys. Rev. D4 (1971) 3497
- [8] D.J. Gross and R. Jackiw, Phys. Rev. D6 (1972) 477
- [9] E. Witten, Nucl. Phys. B223 (1983) 433
- [10] J. Goldstone and F. Wilczek, Phys. Rev. Lett. 47 (1981) 986

BERRY PHASES, MAGNETIC MONOPOLES, AND WESS-ZUMINO TERMS OR HOW THE SKYRMION GOT ITS SPIN*

BY I. J. R. AITCHISON

CERN, Geneva, Switzerland

(Received July 31, 1986)

An elementary discussion is given of the mechanism whereby the Wess-Zumino term determines the quantization of the Skyrme soliton. The work of Balachandran et al. is drawn upon to make explicit the remark of Wu and Zee that the Wess-Zumino term acts like a monopole in the space of scalar fields of the non-linear σ -model. The origin of the monopole structure, and its influence on quantization, is discussed in terms of the Berry (adiabatic) phase.

PACS numbers: 11.17.+y

1. Introduction and outline

The thing about Skyrmions [1] that is surely hardest to understand is how a lump-like solution (soliton) of a classical scalar field theory can, and in some cases even *must*, be quantized as a fermion. How can you add integers together and get a half left over? I want to draw together here some recent work on this subject, which has certainly helped me to understand how this marvellous trick is pulled.

It has of course been known for quite some time that a classical *extended* object (for example, a top [2, 3]) may be quantized as a fermion. A system which provides an explicit model of how this can come about — and one which is directly relevant to Skyrmions — is that of a particle of charge e in motion about a fixed magnetic monopole of strength g . Almost immediately after Dirac's 'monopole' paper [4], Tamm [5] studied the Schrödinger equation for this system, and found that the solutions of the angular equation are the rotation functions $D_{m'm}^j(\theta, \phi)$, with $m = eg$ (in units $\hbar = c = 1$); when the product eg has the minimum non-zero value

$$eg = \frac{1}{2} \quad (1.1)$$

allowed by the Dirac quantization argument [4], or more generally the value $(n + 1/2)$, the system has half-odd angular momentum and is a fermion. (Sometimes this circumstance

* This is an expanded version of the second of two lectures given at the XXVI Cracow School of Theoretical Physics, Zakopane, Poland, 1–13 June, 1986.

is used to run the argument the other way — i.e. that quantization of angular momentum yields the Dirac condition — but, in the present context at least, the Dirac condition will be fundamental.) More recently, field-theory examples of the charge-monopole system have been studied, with analogous results [6, 7].

In the two papers which initiated the recent burst of activity on Skyrmions [1] (and much else besides), Witten [8, 9] showed that the Wess-Zumino ($W - Z$) term [8–10] in the action for the scalar fields ϕ_a (whose solitons are Skyrmions) actually *determines* how these solitons are to be quantized. He obtained the remarkable result that the Skyrmion is a fermion if N_c is odd, and a boson if N_c is even: furthermore, the $W - Z$ term also determines the pattern of spin-SU(3) multiplets ([$1/2^+$, **8**], [$3/2^+$, **10**]...) in the baryon spectrum [9, 11, 12]. Though obviously correct mathematically, these results were nevertheless still hard to explain in physical terms, especially to anyone who did not know what a $W - Z$ term was — and even to those who did¹.

A good deal of light has been shed on this by the work of Balachandran and collaborators [13–17], Berry [18], Stone [19], and Wu and Zee [20, 21]. I shall try to state the major ideas in single sentences, which we will then examine in greater detail in the following sections.

- (i) The $W - Z$ term is a generalization, to the configuration space of scalar fields ϕ_a , of the charge-monopole interaction term in ordinary configuration space for particles. It acts like a monopole in ϕ -space.
- (ii) Because Skyrmion field configurations are maps between field space and real space, the monopole structure of the $W - Z$ term in field space induces, for such configurations, monopole structure in real space.
- (iii) Upon quantization, fermionic behaviour will emerge via the well-known monopole mechanism referred to above.

These sentences state where we are trying to go, but they do not explain (a) where the $W - Z$ term itself comes from, or (b) why it is like a monopole in ϕ -space. The short answer to (a) is: from the very fermion determinant which we studied in the previous lecture, but generalized to SU(3)_f, i.e. it is a term in the effective action for the ϕ fields which arises after integrating over the fermions [22, 23]. This is all very well in its way, but it too is mysterious: why does such an exotic term get induced in the boson sector when we integrate out the fermions? The technical answer to *this* is that the underlying fermion theory has anomalies, which can be calculated from single fermion loop diagrams. These diagrams generate effective vertices in the external fields (ϕ_a , gauge fields, etc.) coupled to the fermions. Hence any bosonic action obtained by integrating out the fermions — which is equivalent to summing all single fermion loop diagrams — *must* faithfully represent these anomaly-induced vertices. The $W - Z$ action precisely encodes these anomalous vertices: if we only consider the ‘ungauged’ $W - Z$ action, which is a function of the SU(3)_f chiral field ϕ alone, we are representing correctly just the SU(3)_f flavour anomalies of the underlying Fermion theory.

¹ For those who know that there is no $W - Z$ term if the flavour group is SU(2), and wonder what happens then, see Section 6.

But anomalies are pretty mysterious too—are we not getting into an infinite regress of ‘explanations’? It would be nice to have some kind of quantum mechanical *analogue*, at least, for what is going on. We can get a clue what to look for when we remember that the characteristic thing about anomaly-induced vertices is that they are independent of the fermion mass M ; it is precisely this circumstance that allows the anomaly-cancellation mechanism discussed in the previous lecture, to work. Thus these ‘anomalous’ vertices will still survive in the fermion determinant with the correct coefficients, even as M becomes very large. This means that these particular vertices—or, equivalently, these particular contributions to the induced bosonic action — can be reliably calculated by the derivative expansion technique: $\partial\phi/M$ can be made as small as we like. (Some explicit examples of this way of calculating anomalous vertices are given in Ref. [22].) Now, a very large fermion mass M implies a large *gap* between the negative energy (sea) levels and the positive energy levels. Small values of $\partial\phi/M$ mean that the momenta and energies associated with these ‘slowly’ varying ϕ fields are much less than the mass gap, and will therefore not induce significant fermionic excitations across the gap — indeed, in the limit of M *very* large, there will be no excitations at all, and we need only deal with the fermion vacuum (ground state).

This state of affairs is something we can find a quantum mechanical (rather than quantum field theoretic) analogue for. It arises quite frequently in many-body physics. Suppose we have a system described in terms of two sets of degrees of freedom: one (which we call r) is ‘fast moving’ with ‘large’ differences between excitation levels, and the other (R) is ‘slow moving’ with ‘small’ associated energy differences. We may think of the electronic (fast) and nuclear (slow) degrees of freedom (d.f.s) in a molecule for instance. It should make sense, when considering the r coordinates, to regard the R ’s as approximately constant; indeed this is called the adiabatic, or Born-Oppenheimer approximation in quantum mechanics. More precisely, if the R were constant, we would simply solve the stationary state Schrödinger equation for the r ’s, with the R ’s appearing parametrically:

$$H_r(R)\psi_n(r, R) = E_n(R)\psi_n(r, R). \quad (1.2)$$

In reality, the R ’s are varying slowly with time, but not quickly enough to induce transitions from one E_n level to another. Thus the system, if started in a particular E_n level, ‘stays with it’ as the R ’s change. This is essentially the content of the adiabatic theorem: the ‘fast’ coordinates stay in the original eigenstate, which however itself changes slowly in response to the slow changes in the R coordinates which appear parametrically. This sounds very much like the situation of our fermion vacuum evolving slowly in response to the slowly varying ϕ ’s. But where is the quantum-mechanical analogue of the W—Z term? It must correspond to some non-trivial structure left behind in the space of the ‘slow’ parameters when we adiabatically decouple the ‘fast’ ones.

Here is where the work of Berry [18], and Kuratsuji and Iida [24] comes in. The adiabatic assumption tells us that, at any time t , the state of the system $|\psi(t)\rangle$ (adopting now a slightly more abstract notation) will essentially be the ‘instantaneous’ eigenstate $|n(R(t))\rangle$, where

$$H(R(t))|n(R(t))\rangle = E_n(R(t))|n(R(t))\rangle, \quad (1.3)$$

210

if it was prepared to be in one of these states $|n(\mathbf{R}_0)\rangle$ at $t = 0$, where $\mathbf{R}_0 = \mathbf{R}(t = 0)$. In fact, $|\psi(t)\rangle$ will be related to $|n(\mathbf{R}(t))\rangle$ by a phase factor. What phase factor? The naïve answer would surely be

$$|\psi(t)\rangle = [\exp -i \int_0^t E_n(\mathbf{R}(t')) dt'] \cdot |n(\mathbf{R}(t))\rangle, \quad (1.4)$$

the expected integrated ‘quasi-stationary state’ phase. But this is *not* the whole story. An *additional* phase is generated during such an adiabatic change. That this is so in principle has been known for a long time (see, for example, Ref. [25]), but it had tended to be dismissed as unimportant physically (‘just a phase’). Berry [18] pointed out a number of cases where the phase could be of considerable physical interest (see also Ref. [26]). In particular, a non-trivial phase can result from a closed path in \mathbf{R} space, as we move along $\mathbf{R}_0 \rightarrow \mathbf{R}(t) \rightarrow \mathbf{R}_0$. Such ‘Berry phases’ depend on the actual path followed in \mathbf{R} -space — which may remind us of something...

The Berry phase is easily calculated [18]. We are looking for a solution of

$$H(\mathbf{R}(t)) |\psi(t)\rangle = i \frac{d}{dt} |\psi(t)\rangle, \quad (1.5)$$

and we try the adiabatically-inspired ansatz

$$|\psi(t)\rangle = [\exp -i \int_0^t E_n(\mathbf{R}(t')) dt'] \cdot \exp i\gamma_n(t) \cdot |n(\mathbf{R}(t))\rangle. \quad (1.6)$$

Inserting (1.6) directly into (1.5) and using (1.3) yields

$$\begin{aligned} \gamma_n(t) &= i \int_0^t \langle n(\mathbf{R}(t')) | \frac{d}{dt'} |n(\mathbf{R}(t'))\rangle dt' \\ &= i \int_0^t \langle n(\mathbf{R}) | \nabla_{\mathbf{R}} n(\mathbf{R}) \rangle \cdot d\mathbf{R}, \end{aligned} \quad (1.7)$$

the fundamental formula [18] for the Berry phase $\gamma_n(t)$.

Now — having dealt adiabatically with the \mathbf{r} d.f.’s this way --- let us turn our attention to the slowly varying \mathbf{R} ’s, and consider *them* as quantum d.f.’s, not just classical parameters. In the same adiabatic approximation, we sit in one ‘electronic’ state n and look for solutions in which the total state function has the product form $\phi_n(\mathbf{R})|n(\mathbf{R})\rangle$, and ask: what Schrödinger equation does $\phi_n(\mathbf{R})$ obey? The answer is very interesting [25]: if $V(\mathbf{R})$ is the potential energy relevant to the \mathbf{R} d.f.’s alone, then $\phi_n(\mathbf{R})$ obeys

$$[\text{‘covariant kinetic energy’} + E_n(\mathbf{R}) + V(\mathbf{R})] \phi_n(\mathbf{R}) = i \frac{d}{dt} \phi_n(\mathbf{R}), \quad (1.8)$$

where by ‘covariant kinetic energy’ is meant that the gradient operator $\nabla_{\mathbf{R}}$ in the normal \mathbf{R} -kinetic energy terms is replaced by

$$\nabla_{\mathbf{R}} \rightarrow \nabla_{\mathbf{R}} + \langle n(\mathbf{R}) | \nabla_{\mathbf{R}} n(\mathbf{R}) \rangle \quad (1.9)$$

$$\equiv \nabla_{\mathbf{R}} - iA_n(\mathbf{R}), \quad (1.10)$$

where (1.10) follows since the matrix element in (1.8) is easily seen to be pure imaginary. Thus a sort of gauge potential has been induced in \mathbf{R} -space!

It is clear that this gauge potential is intimately related to the Berry phase; they are two facets of the same subtlety in the adiabatic approximation. Indeed we can see from (1.9) and (1.10) exactly what the local phase invariance corresponding to this ‘gauge’ structure is: an \mathbf{R} -dependent phase change on $|n(\mathbf{R})\rangle$ induces a change in A_n of (1.10), which in turn causes a precisely compensating phase change in $\phi_n(\mathbf{R})$, so that the total state function $\phi_n(\mathbf{R})|n(\mathbf{R})\rangle$ is locally phase invariant. Thus a distinctly non-trivial structure has appeared in the ‘slow’ d.f.’s after adiabatic decoupling of the ‘fast’ d.f.’s. Note, incidentally, that the Berry phase is nothing but

$$\exp i\gamma_n(t) = \exp [i \int_0^t A_n(\mathbf{R}) \cdot d\mathbf{R}], \quad (1.11)$$

so that we were right to be reminded of the *path-dependence* of the wave function for a particle in an electromagnetic potential A .

Thus we are getting nearer to understanding how funny phase factors — which might influence apparent rotational properties [26] — can arise via adiabatic decoupling. We can make closer contact with the field theory if we reformulate the adiabatic approximation in the path integral formalism. This was done by Kuratsuji and Iida [24]. In view of (1.6) and (1.11) we can almost guess what the result must be. We want the dynamics in \mathbf{R} -space to correspond to a ‘particle’ moving in an (additional) ‘vector potential’ $A_n(\mathbf{R})$. Thus we expect to find a piece in the effective action $S_{\text{eff},n}$ in \mathbf{R} -space which corresponds to the effective Lagrangian

$$\mathcal{L}_{\text{eff},n} = A_n(\mathbf{R}) \cdot \frac{d\mathbf{R}}{dt}. \quad (1.12)$$

Indeed, in that case

$$S_{\text{eff},n}(T) = \int_0^T \mathcal{L}_{\text{eff},n} dt = \int_0^T A_n(\mathbf{R}) \cdot \frac{d\mathbf{R}}{dt} dt, \quad (1.13)$$

and

$$\exp iS_{\text{eff},n}(T) = \exp i\gamma_n(T). \quad (1.14)$$

This is just what Kuratsuji and Iida obtain. By considering the trace of the evolution operator $\text{tr exp}(-iHT)$ in the adiabatic approximation, they show that it is given by

$$K_{\text{eff}}(T) = \sum_n \int \mathcal{D}\mathbf{R} \exp \{iS_0 - i \int_0^T E_n(\mathbf{R}) dt' + i\gamma_n(T)\}, \quad (1.15)$$

where $\mathbf{R}(T) = \mathbf{R}_0$ (since for the trace we want to return to the same state at $t = T$), and where γ_n is now evaluated over closed loops $\mathbf{R}_0 \rightarrow \mathbf{R}(t) \rightarrow \mathbf{R}(T) = \mathbf{R}_0$ in \mathbf{R} -space:

$$\gamma_n(T) = i \oint \langle n(\mathbf{R}) |\nabla_{\mathbf{R}} n(\mathbf{R}) \rangle \cdot d\mathbf{R} = \oint \mathbf{A}_n(\mathbf{R}) \cdot d\mathbf{R} = \oint \mathcal{L}_{\text{eff},n} dt. \quad (1.16)$$

S_0 is the ordinary action for the \mathbf{R} coordinates.

Now, finally, how can we understand the *specific* ‘monopole-like’ structure which corresponds (we have asserted) to the $W - Z$ term? The secret, as Stone [19] pointed out, lies in a beautiful discovery by Berry [18]. We have assumed throughout that the eigenvalues $E_n(\mathbf{R})$ were well separated, and certainly not degenerate. But what happens if, for some particular value of the \mathbf{R} d.f.’s, say \mathbf{R}^* , two of the E_n ’s coalesce? We expect some sort of catastrophe to show up in our adiabatic result. In fact, this point in \mathbf{R} -space is very likely to be a point at which the vector potential $\mathbf{A}_n(\mathbf{R})$ is singular! Such a vector potential would imply sources (δ -function singularities in the associated field strengths) — for example, magnetic monopoles. This is exactly what Berry found, explicitly, for the case in which the degeneracy is a spin-type degeneracy, the ‘fast’ coordinates are spin d.f.’s, and the ‘slow’ ones are angles describing the orientation of the (real!) magnetic field \mathbf{B} . The equation corresponding to (1.5) is then

$$\mu \mathbf{BS} \cdot \hat{\mathbf{B}} |\psi(t)\rangle = i \frac{d}{dt} |\psi(t)\rangle \quad (1.17)$$

and

$$E_n(\mathbf{B}) = \frac{1}{2} \mu B_n, \quad (1.18)$$

where $n/2$ is the spin eigenvalue, which takes $2s+1$ values. Clearly these $2s+1$ states are degenerate when $\mathbf{B} = \mathbf{0}$ (the point \mathbf{R}^*). Berry found that the associated $\mathbf{A}_n(\mathbf{B})$ was precisely that of a monopole in \mathbf{B} -space located at $\mathbf{B} = \mathbf{0}$, having strength $-n/2$ (i.e. $eg = -n/2$). Thus spin-type degeneracies cause monopoles to lurk in the ‘slow’ space.

We can now see why the integral in (1.16) along a closed loop need not vanish. If we convert the line integral in (1.16) by a (multi-dimensional) Stokes theorem to a surface integral over the ‘magnetic field’ $\mathbf{B}_n = \nabla \times \mathbf{A}_n$, and thence to a volume integral via Gauss, we would normally get zero since $\text{div } \mathbf{B} = 0$. However, for the singular potential corresponding to a monopole $\text{div } \mathbf{B}_n \neq 0$ and a closed loop contributes a non-zero result. Actually we can go even further than this. The line integral over a closed loop C becomes

$$\oint_C \mathbf{A}_n(\mathbf{R}) \cdot d\mathbf{R} = \iint_S \mathbf{B}_n \cdot d\mathbf{S}, \quad (1.19)$$

where S is a surface spanning C . But what surface? Should we take an S_1 (see Fig. 1) which is ‘above’ C , or an S_2 which is ‘below’? For consistency we must have

$$\iint_{S_1} \mathbf{B}_n \cdot d\mathbf{S} = \iint_{S_2} \mathbf{B}_n \cdot d\mathbf{S} + 2N\pi \quad (1.20)$$

(remember that these quantities are all *phases*). Since the normals for S_2 and for S_1 are oppositely oriented, we see that (1.20) is equivalent to

$$\oint \mathbf{B}_n \cdot d\mathbf{S} = 2N\pi, \quad (1.21)$$

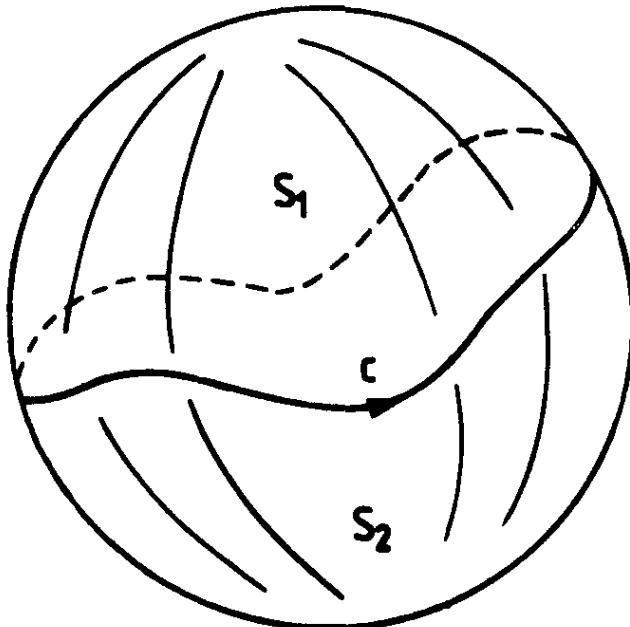


Fig. 1. Two surfaces spanning the curve C on S_2

where the integral is over a *closed* surface surrounding \mathbf{R} . Thus the total flux out of the ‘monopole’ is quantized — which is just the Dirac [3] condition (*eg* = $N \cdot 1/2$, in the real electromagnetic (e.m.) case). The above argument was a poor man’s version of a deeper topological treatment, since we relied heavily on tacitly thinking of \mathbf{R} as three-dimensional. Nevertheless, the result is correct.

But now notice a remarkable thing: the previous paragraph has shown quite generally that the monopole strength (total flux through a closed surface, divided by 4π) has the quantized value $N/2$, where N is an integer of topological significance. The paragraph before that stated the result that monopole-like structure arose from (1.17), in which the strength of the monopole is $-n/2$, where $n/2$ is the spin eigenvalue. The eigenvalue spectrum of (1.17) near $\mathbf{B} = \mathbf{0}$ (the point of degeneracy) seems to know something about topology!

We have learned that monopole-like structure can be generated in the ‘slow’ space \mathbf{R} when we adiabatically decouple the ‘fast’ d.f.’s. Furthermore, the strength of the monopole interaction is an integer (divided by 2), from topological considerations, and this integer corresponds to some label of the energy spectrum of the ‘fast’ coordinates. This is as near as we are likely to get to a quantum mechanical analogue of the W-Z mystery. The W-Z term results [22, 23] from adiabatically decoupling the ψ ’s from the π ’s, starting from a Dirac equation

$$[-i\alpha \cdot \nabla + \beta\mu \exp(i\lambda \cdot \pi\gamma_5 f^{-1})]\psi = i \frac{\partial\psi}{\partial t}, \quad (1.22)$$

which is the analogue of (1.17); μ is a mass parameter, $f \approx 93$ MeV, λ_a ($a = 1, 2, \dots, 8$) are the Gell-Mann matrices, and the eight π fields are the analogues of the angle variables in $\hat{\mathbf{B}}$. The W-Z term (in the fields π) looks like a monopole in π space. Its coefficient is found to be an integer (which is, of course, N_c) by topological considerations [8] exactly

analogous to those given above for the \mathbf{R} -space monopole. There is one gap left to be closed: what is it in the spectrum of the Dirac equation (1.22) that ‘knows’ about topology (and hence about monopoles)? That is a *deep* question, the answer to which is provided by the mathematical subject called index theory. This way of looking at anomalies (remember?) is called the ‘Hamiltonian approach’ [27, 28], and is precisely the quantum field theoretical analogue of the quantum mechanical Berry-phase discussion outlined above.

Let us now see how all the foregoing works out in some simple cases.

2. A simple example

Consider, following Stone [19], a spin-1/2 particle in a magnetic field $\mathbf{B} = B\mathbf{n}$, where $\mathbf{n}^2 = 1$. The state function for the (‘fast’) spin d.f.’s satisfies

$$\mu\boldsymbol{\sigma} \cdot \mathbf{n}(t)|\psi(t)\rangle = i \frac{d}{dt} |\psi(t)\rangle, \quad (2.1)$$

where the magnitude B of the magnetic field has been absorbed into μ . The ‘slow’ d.f.’s are \mathbf{n} , since we shall only consider slow variations of \mathbf{B} with fixed B . We consider the large μ limit (cf. large μ in (1.22)), so that the slow changes in \mathbf{n} do not cause transitions between the two spin eigenstates $|\uparrow\mathbf{n}(t)\rangle$ and $|\downarrow\mathbf{n}(t)\rangle$, in the adiabatic approximation. Suppose at $t = 0$ we start in the state $|\uparrow\mathbf{n}(0)\rangle$, where $\mathbf{n}(0) = (\sin\theta_0 \cos\phi_0, \sin\theta_0 \sin\phi_0, \cos\theta_0)$. At a general time t , $\mathbf{n}(t) = (\sin\theta \cos\phi, \sin\theta \sin\phi, \cos\theta)$, where the time-dependent d.f.’s θ and ϕ will vary over the surface of an S_2 . The Berry phase $\gamma_f(t)$ is

$$\gamma_f(t) = i \int_0^t \langle \mathbf{n}(t') \uparrow \left| \frac{d}{dt'} \right| \uparrow \mathbf{n}(t') \rangle dt'. \quad (2.2)$$

We shall calculate this directly using an explicit wave function for $|\uparrow\mathbf{n}\rangle$; already here a crucial feature will emerge.

The wave function

$$\langle \theta\phi | \uparrow\mathbf{n}+ \rangle = \begin{pmatrix} \cos\theta/2 \\ \sin\theta/2 e^{i\phi} \end{pmatrix} \quad (2.3)$$

is certainly an eigenfunction of

$$\mu\boldsymbol{\sigma} \cdot \mathbf{n} = \mu \begin{pmatrix} \cos\theta & \sin\theta e^{-i\phi} \\ \sin\theta e^{i\phi} & -\cos\theta \end{pmatrix} \quad (2.4)$$

with eigenvalue μ . Inserting (2.3) into (2.2) we find

$$\gamma_f^+(t) = - \int_0^t \frac{1}{2} (1 - \cos\theta) \frac{d\phi}{dt'} dt' \quad (2.5)$$

for the Berry phase $\gamma_{\uparrow}^+(t)$, as θ and ϕ vary slowly over S_2 . The reason for the + symbols will become clear in a moment.

According to what was advertised in Section 1, the integrand in (2.5) should be closely related to the vector potential of a magnetic monopole of strength $-1/2$, in θ - ϕ space, positioned at the origin. A standard expression for the vector potential of a Dirac monopole of this strength is

$$\mathbf{A}_+(\mathbf{r}) = \frac{-1}{2r} \frac{1}{z+r} \cdot (-y, x, 0), \quad (2.6)$$

where $x = r \sin \theta \cos \phi$, $y = r \sin \theta \sin \phi$, $z = r \cos \theta$. Thus with $\mathbf{r} = (x, y, z)$

$$\mathbf{A}_+(\mathbf{r}) \cdot d\mathbf{r} = \frac{-1}{2r(z+r)} (xdy - ydx) = -\frac{1}{2}(1 - \cos \theta)d\phi. \quad (2.7)$$

Hence indeed (cf. (1.11))

$$\gamma_{\uparrow}^+(t) = \int_0^t \mathbf{A}_+(\mathbf{n}) \cdot \frac{d\mathbf{n}}{dt'} dt' \equiv \int_0^t \mathcal{L}_{\text{eff}}^+(\mathbf{n}) dt', \quad (2.8)$$

and we see explicitly the ‘monopole’ character of the phase factor associated with adiabatic motion round the degeneracy point $\mathbf{n} = \mathbf{0}$. (We hope the reader will not be confused by the use of \mathbf{n} for the slow d.f.’s in this Section, and of n as a label of a ‘fast’ eigenstate in the previous one.)

However, the wave function (2.3) is ill-defined at $\theta = \pi$ (what is the value of ϕ when $\theta = \pi$?). An alternative choice of \uparrow wave function which is well defined at $\theta = \pi$ is

$$\langle \theta\phi | \uparrow \mathbf{n} - \rangle = \begin{pmatrix} \cos \theta/2 e^{-i\phi} \\ \sin \theta/2 \end{pmatrix}. \quad (2.9)$$

Repeating the above calculations we find that this leads to a Berry phase

$$\gamma_{\uparrow}^-(t) = - \int_0^t \frac{1}{2} (-1 - \cos \theta) \frac{d\phi}{dt'} dt', \quad (2.10)$$

which is equivalent to a vector potential

$$\mathbf{A}_-(\mathbf{r}) = \frac{1}{2r} \cdot \frac{1}{z-r} (-y, x, 0). \quad (2.11)$$

Though good at $\theta = \pi$, (2.9) is ill-defined at $\theta = 0$ — and in fact we are hitting here the famous problem that, for a monopole field, no *single* vector potential exists which is singularity-free over the entire manifold S_2 . The \mathbf{A}_+ which followed from the choice (2.3) has

a singularity along $z = -r$, i.e. the negative z -axis, or $\theta = \pi$. This line of singularities is called a ‘Dirac string’ [4]. Likewise, the A_+ choice has a string along $\theta = 0$. But, comparing (2.5) and (2.10) we see that A_+ and A_- differ by a gradient

$$A_+ - A_- = \nabla\phi. \quad (2.12)$$

that is, by a gauge transformation. Correspondingly,

$$\gamma_t^+ - \gamma_t^- = -[\phi(t) - \phi(0)], \quad (2.13)$$

so that for a closed path on S_2 , γ_t is unique.

What we see explicitly here for A_\pm is generally true. Any particular A will have a string singularity somewhere, and by doing a gauge transformation we merely shift the singularity somewhere else. The use of *two* A ’s (e.g. A_+ and A_-) was advocated by Wu and Yang [29, 30] as a way round the singularity problem, since we can use each in a region (or ‘patch’) where it is singularity-free, and then connect the two, in a convenient overlap region, by a gauge transformation.

The problem of singularities in the vector potential corresponding to a magnetic monopole would seem to be unavoidable since, if $B = \nabla \times A$ and A is singularity-free, $\text{div } B = 0$ and the magnetic charge must be zero. In our ultimate application of the Berry phase concept, the ‘slow’ d.f.’s will be the meson field variables, which we shall want to quantize. This is analogous to quantizing the n d.f.’s (i.e. θ, ϕ) in the present quantum-mechanical analogue. The presence of the (monopole) singularity at $n = 0$ makes this quantization very awkward, and the Wu–Yang procedure is also not well-adapted to our later purpose.

Remarkably enough, however, it is possible to find a singularity-free Lagrangian for the monopole problem. Indeed, it is given to us automatically by the Berry phase formula, as we shall now describe. We then show how to obtain, from the Berry formula, the elegant Balachandran formalism [13–17], which is ideally suited to the Skyrmiion application.

3. The Hopf fibration of S_2 , and the Balachandran Lagrangian

Let us introduce the notation

$$z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \quad (3.1)$$

for the two-component spinor which is the eigenfunction of (2.1) (e.g. z could be (2.3) or (2.9)). The effective Lagrangian associated with the Berry phase is then

$$\mathcal{L}_{\text{eff}} = iz^\dagger \frac{dz}{dt}. \quad (3.2)$$

(cf. (2.8) and (2.10), and check it by trying (2.3) or (2.9) for z). In the two z ’s considered explicitly so far ((2.3) and (2.9)) only two d.f.’s entered, namely θ and ϕ , the coordinates of a point on the surface of a two-dimensional sphere S_2 . We set $\mathcal{L}_{\text{eff}} = -A \cdot dn/dt$ to

obtain the potentials A_+ and A_- , also on S_2 . However, in principle the spinor z has *three* d.f.'s, since the normalization condition

$$z^\dagger z = 1 = |z_1|^2 + |z_2|^2 \quad (3.3)$$

is only one constraint on the two complex numbers z_1, z_2 . Indeed, we may in general consider either (2.3) or (2.9) to be multiplied by an arbitrary *phase*, for example

$$z = \begin{pmatrix} \cos \theta/2 & e^{i\chi} \\ \sin \theta/2 & e^{i(\phi+\chi)} \end{pmatrix}. \quad (3.4)$$

The corresponding ' A ' must now depend on three d.f.'s, and consequently is not restricted to the surface of an S_2 : it turns out, as we shall now see, that it is actually defined on the surface of an S_3 , and is non-singular!

Suppose we write

$$\left. \begin{aligned} z_1 &= x_1 + ix_2 \\ z_2 &= x_3 + ix_4 \end{aligned} \right\}. \quad (3.5)$$

Then $z^\dagger z = 1$ becomes

$$x_1^2 + x_2^2 + x_3^2 + x_4^2 = 1, \quad (3.6)$$

and the x_i are the coordinates of a point on an S_3 . Comparing (3.4) and (3.5) we find

$$\left. \begin{aligned} x_1 &= \cos \theta/2 \cos \chi, & x_2 &= \cos \theta/2 \sin \chi \\ x_3 &= \sin \theta/2 \cos (\phi + \chi), & x_4 &= \sin \theta/2 \sin (\phi + \chi) \end{aligned} \right\}. \quad (3.7)$$

The metric is

$$ds^2 = \frac{1}{4} d\theta^2 + d\chi^2 + \sin^2 \theta/2 d\phi^2 + 2 \sin^2 \theta/2 d\phi d\chi. \quad (3.8)$$

It is more convenient to use orthogonal coordinates by introducing

$$\psi = \phi + \chi \quad (3.9)$$

in terms of which (3.8) becomes

$$ds^2 = \frac{1}{4} d\theta^2 + \cos^2 \theta/2 d\chi^2 + \sin^2 \theta/2 d\psi^2. \quad (3.10)$$

Thus ' $A \cdot dn$ ' now has the form

$$A_\theta \frac{1}{2} d\theta + A_\chi \cos \theta/2 d\chi + A_\psi \sin \theta/2 d\psi. \quad (3.11)$$

Inserting (3.4) into (3.2) we find easily,

$$'A \cdot dn' = -iz^\dagger dz = d\chi + \frac{1}{2}(1 - \cos \theta)d\phi \quad (3.12)$$

$$= \cos^2 \theta/2 d\chi + \sin^2 \theta/2 d\psi \quad (3.13)$$

whence, via (3.11),

$$A_\theta = 0, \quad A_\chi = \cos \theta/2, \quad A_\psi = \sin \theta/2. \quad (3.14)$$

These potentials are manifestly non-singular. By contrast, the ' S_2 ' forms (2.7), and the corresponding $\mathbf{A}_- \cdot d\mathbf{n}$ from (2.11), are singular. Consider, for example (2.7). On S_2 the metric is $ds^2 = d\theta^2 + \sin^2 \theta d\phi^2$ and so

$$A_{+, \theta} = 0, A_{+, \phi} = -\frac{1}{2} \frac{(1 - \cos \theta)}{\sin \theta} = -\frac{1}{2} \tan \theta/2, \quad (3.15)$$

which is singular (as expected) at $\theta = \pi$. Likewise $A_{-, \phi}$ is singular at $\theta = 0$. In terms of the S_3 coordinates,

$$\mathbf{A}_+ \cdot d\mathbf{n} = -\frac{1}{2} (1 - \cos \theta) d\phi = -\sin^2 \theta/2 d\psi + \sin^2 \theta/2 d\chi, \quad (3.16)$$

giving

$$A_{+, \theta} = 0, \quad A_{+, \chi} = \frac{\sin^2 \theta/2}{\cos \theta/2}, \quad A_{+, \psi} = -\sin \theta/2 \quad (3.17)$$

and $A_{+, \chi}$ is singular at $\theta = \pi$. The S_3 components of \mathbf{A}_- can be found similarly, and this time $A_{-, \psi}$ is singular at $\theta = 0$.

Thus a non-singular potential for the monopole can be found provided we enlarge the configuration space from S_2 to S_3 , and use the full three d.f.'s available in z . Are we sure that the physics is really the same? The Lagrangian \mathcal{L}_{eff} corresponding to (3.4) is, of course,

$$\mathcal{L}_{\text{eff}} = \mathcal{L}_{\text{eff}}^+(\mathbf{n}) - \dot{\chi}, \quad (3.18)$$

which differs from $\mathcal{L}_{\text{eff}}^+(\mathbf{n})$ by a total time derivative, and therefore leads to the same equations of motion. From (3.18) we learn that χ is acting like a U(1) gauge d.f. Thus the two d.f.'s of S_2 have been enlarged to three by the addition of a U(1) gauge d.f., χ . This is a well-known construction in mathematics, called the Hopf fibration of S_2 . S_3 can be regarded as a principal fibre bundle with base space S_2 and a U(1) structure group. The (Hopf) projection map which takes us from S_3 to S_2 is given explicitly by

$$\mathbf{n} = z^\dagger \sigma z, \quad (3.19)$$

as can easily be checked (Appendix B). Ryder [31] and Minami [32] were the first to introduce the Hopf map into monopole theory.

From (3.12) it is clear that the \mathbf{A}_+ potential is obtained (cf. (2.7)) by setting $\chi = 0$, and the \mathbf{A}_- one by setting $\chi = -\phi$. Restricting χ in this way is called taking a 'section' of the fibre bundle. These two choices are each called 'local' sections, because they (and the potentials) are not smoothly defined globally over the entire S_2 : \mathbf{A}_+ is smooth for an upper patch of S_2 excluding $\theta = \pi$, and \mathbf{A}_- is smooth for a lower patch excluding $\theta = 0$. It is, in fact, not possible to find any such section which is smooth globally, in this case: a minimum of two is required, as in the explicit examples of \mathbf{A}_\pm . Mathematically this corresponds to the fact that our (monopole) bundle is non-trivial, or — equivalently — to the fact that S_3 is only locally, but not globally, equivalent to $S_2 \times S_1$. Thus the monopole Lagrangian can be described in a singularity-free way by using a non-trivial bundle over S_2 .

The above formulation is not yet quite suitable for our later application to Skyrmion physics. In that case, the d.f.'s in which we shall be interested are actually entries in an $SU(3)$ matrix, and it is hard to see how to generalize z to such a matrix. On the other hand, S_3 is the group manifold of $SU(2)$, and it is quite simple to reformulate the above results in terms of a basic dynamical variable $s(\theta, \phi, \chi) \in SU(2)$, rather than $z(\theta, \phi, \chi)$. This will lead to Balachandran's form for \mathcal{L}_{eff} , which will be directly analogous to the $SU(3)$ case.

We can associate a general $SU(2)$ matrix s with the components z_1, z_2 of z via

$$s = \begin{pmatrix} z_1 & -z_2^* \\ z_2 & z_1^* \end{pmatrix} \quad (3.20)$$

since the condition $|z_1|^2 + |z_2|^2 = 1$ guarantees $s^\dagger s = ss^\dagger = 1$. In terms of s , the \mathcal{L}_{eff} of (3.2) becomes

$$\mathcal{L}_{\text{eff}} = iz^\dagger \frac{dz}{dt} = \frac{i}{2} \text{tr} (\sigma_3 s^{-1} \dot{s}) \quad (3.21)$$

as may be verified explicitly. Equation (3.21) provides our desired (Balachandran) monopole Lagrangian in terms of $s \in SU(2)$. It is pleasing to see this direct link between the Berry phase and the Balachandran Lagrangian.

The Hopf map can equivalently be described in terms of s . The counterpart of (3.19) is

$$\sigma \cdot n = s \sigma_3 s^{-1} \quad (3.22)$$

(note that $n^2 = 1$ follows automatically upon squaring both sides). Under right multiplication of s by an element of $U(1)$

$$s \rightarrow s \exp i\sigma_3 \alpha, \quad (3.23)$$

$$\sigma \cdot n \rightarrow s(\exp i\sigma_3 \alpha) \sigma_3 (\exp -i\sigma_3 \alpha) s^{-1} = s \sigma_3 s^{-1} = \sigma \cdot n \quad (3.24)$$

and n is unchanged. Thus the space $SU(2)/U(1)$ of right cosets (3.23) gets mapped by (3.22) into S_2 (see Fig. 2).

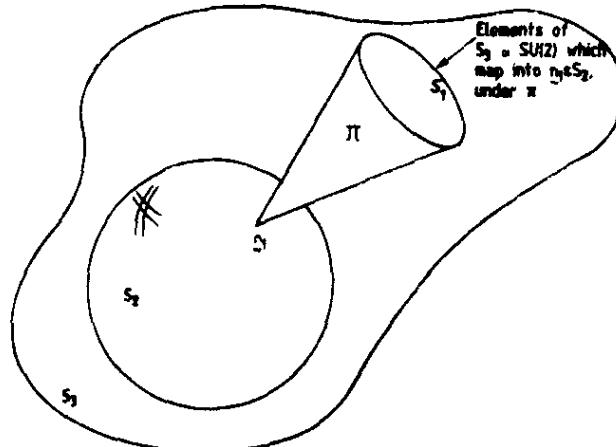


Fig. 2. The one-cycle S_1 of all points in S_3 related to a given point s_1 of S_3 by $s_1 \exp(i\sigma_3 \alpha)$, as α varies, is mapped by the Hopf map π into the single point n_1 of S_2 .

To gain some confidence with this s -formalism, we can simply insert (2.3) into (3.20), obtaining

$$s_+(\mathbf{n}) = \begin{pmatrix} \cos \theta/2 & -\sin \theta/2 e^{-i\phi} \\ \sin \theta/2 e^{i\phi} & \cos \theta/2 \end{pmatrix} \quad (3.25)$$

whence

$$\mathcal{L}_{\text{eff}}^+(\mathbf{n}) = \frac{i}{2} \text{tr}(\sigma_3 s_+^{-1} \dot{s}_+) = \frac{-(1-\cos \theta)}{2} \frac{d\phi}{dt} \quad (3.26)$$

in agreement with (2.8). Alternatively, (2.9) gives

$$s_-(\mathbf{n}) = \begin{pmatrix} \cos \theta/2 e^{-i\phi} & -\sin \theta/2 \\ \sin \theta/2 & \cos \theta/2 e^{i\phi} \end{pmatrix} \quad (3.27)$$

so that

$$s_+(\mathbf{n}) = s_-(\mathbf{n}) e^{i\sigma_3 \phi} \quad (3.28a)$$

$$s_+ \sigma_3 s_+^{-1} = s_- \sigma_3 s_-^{-1} = \boldsymbol{\sigma} \cdot \mathbf{n} \quad (3.28b)$$

$$\mathcal{L}_{\text{eff}}^-(\mathbf{n}) = \frac{i}{2} \text{tr}(\sigma_3 s_-^{-1} \dot{s}_-) = \frac{(-1-\cos \theta)}{2} \frac{d\phi}{dt} \quad (3.28c)$$

$$\mathcal{L}_{\text{eff}}^+(\mathbf{n}) - \mathcal{L}_{\text{eff}}^-(\mathbf{n}) = -\dot{\phi}. \quad (3.28d)$$

Equation (3.28a) shows that, for given \mathbf{n} , s_+ and s_- are in the same coset, and get mapped into the same \mathbf{n} (3.28b). Equation (3.28d) shows that the corresponding effective Lagrangians differ by a total derivative, and hence lead to the same equations of motion for the \mathbf{n} d.f.'s. Indeed, the difference between the choice s_+ and s_- corresponds exactly to the gauge transformation on the associated vector potentials \mathbf{A}_+ and \mathbf{A}_- considered earlier in (2.12).

In general, we may consider now the SU(2) matrix

$$\begin{aligned} s(\theta, \phi, \chi) &= s_+(\theta, \phi) e^{i\sigma_3 \chi} \\ &= \begin{pmatrix} \cos \theta/2 e^{i\chi} & -\sin \theta/2 e^{-i(\phi+\chi)} \\ \sin \theta/2 e^{i(\phi+\chi)} & \cos \theta/2 e^{-i\chi} \end{pmatrix}, \end{aligned} \quad (3.29)$$

corresponding to (3.4). Then s_+ is the $\chi = 0$ section of this, while s_- is the $\chi = -\phi$ section. The Lagrangian following from (3.29) is, of course,

$$\frac{i}{2} \text{tr}(\sigma_3 s_+^{-1} \dot{s}_+) - \dot{\chi} \quad (3.30)$$

as in (3.18).

In concluding this Section we note (see [31] and [32]) that the foregoing can all be rephrased using the compact formalism of differential forms, and some elementary ideas

of homology and cohomology. The potential 1-form is

$$A = iz^\dagger dz = -x_1 dx_2 + x_2 dx_1 - x_3 dx_4 + x_4 dx_3, \quad (3.31)$$

and the field 2-form B is

$$B = dA = idz^\dagger \wedge dz = -2(dx_1 \wedge dx_2 + dx_3 \wedge dx_4) = -\frac{1}{2} \sin \theta d\theta \wedge d\phi, \quad (3.32)$$

where (3.7) has been used. B is just proportional to the area 2-form of S_2 , and

$$\int_{S_2} B = -\frac{1}{2} 4\pi, \quad (3.33)$$

showing that these potentials and fields indeed correspond to a monopole of strength 1/2. B is certainly closed,

$$dB = 0$$

but it cannot be exact ($B = dA$) on S_2 , since if it were we could use Stokes' theorem on S_2 to obtain

$$\int_{S_2} B = \int_{S_2} dA = \int_{\partial S_2} A = 0, \quad (3.34)$$

since S_2 has no boundary; (3.34) would then contradict (3.33). However, if B is regarded as a 2-form on S_3 it is exact, since $H^2(S_3) = 0$, and consequently an A such that $B = dA$ does exist on S_3 .

4. Quantization of the n d.f.'s: the Dirac condition again

We now want to consider, following Balachandran et al. [13, 17], the problem of quantizing the d.f.'s $n(\theta, \phi)$ — i.e. we want to promote the ‘slow’ d.f.’s, which hitherto in Sections 2 and 3 have been parameters, to dynamical variables. In path integral terms, this means — cf. (1.14) — that we want to consider

$$\int \mathcal{D}n \exp \{i \int [\frac{1}{2} \dot{n}^2 + \mathcal{L}_{\text{eff}}(n)] dt\}. \quad (4.1)$$

We know that the quantum theory of the charge-monopole system should only be consistent provided the Dirac condition holds. We are going to see where this arises in the s -formalism.

In the previous Section we have seen that the introduction of the new (SU(2)) d.f. χ allowed us to describe the monopole system by a non-singular Lagrangian — and so in quantizing this system we do not have the problem of singularities to contend with. On the other hand, we want the physics to be independent of χ . In the classical theory, as we have seen χ acts like a U(1) gauge d.f., and changing χ is like doing a gauge transformation, under which the equations of motion are invariant. In the quantum theory, we must ensure that a corresponding gauge invariance is correctly implemented. This requirement leads to the Dirac condition.

It is clear that the first term, $1/2\mathbf{n}^2$, in the Lagrangian of (4.1) is *invariant* under a U(1) gauge transformation

$$s \rightarrow se^{i\sigma_3\alpha(t)} \quad (4.2)$$

since \mathbf{n} remains invariant under (4.2) (see also Appendix C). Thus non-trivial constraints on the theory, associated with the implementation of gauge invariance under (4.2), must arise from the second ('monopole') term. Let us consider a general such term

$$\mathcal{L}_{\text{eff}}(\mathbf{n}) = -gi \operatorname{tr}(\sigma_3 s^{-1} \dot{s}), \quad (4.3)$$

where the monopole strength g is not yet determined. Then, under (4.2),

$$\mathcal{L}_{\text{eff}} \rightarrow \mathcal{L}_{\text{eff}} + 2g\dot{\alpha}. \quad (4.4)$$

In the quantum theory, s will be promoted to a quantum variable \hat{s} , and wave functions will be written as $\Psi(s)$. Consider the infinitesimal (quantum) version of (4.2):

$$\hat{s} \rightarrow \hat{s} + i\hat{s}\sigma_3\delta\alpha, \quad (4.5)$$

and let \hat{G} be the generator of this transformation so that

$$[\hat{G}, \hat{s}] = \hat{s}\sigma_3. \quad (4.6)$$

Then, from Noether's theorem and (4.4), we deduce

$$\hat{G}\Psi = 2g\Psi \quad (4.7)$$

as a consistency condition on the state functions (it is a kind of 'Gauss Law' associated with gauge invariance under (4.2); see also Appendix C). For finite transformations we then have

$$\Psi'(s) \equiv (e^{i\hat{G}\alpha}\Psi)(s) = \Psi(se^{i\sigma_3\alpha}) = e^{2ig\alpha}\Psi. \quad (4.8)$$

The last two equalities of (4.8) give

$$\Psi(\theta, \phi, \chi + \alpha) = e^{2ig\alpha}\Psi(\theta, \phi, \chi), \quad (4.9)$$

which enforces a kind of 'Bloch' condition on the χ d.f. If we consider the particular case $\alpha = 2\pi$, then since

$$e^{2\pi i\sigma_3} = 1 \quad (4.10)$$

we deduce

$$e^{4\pi ig} = 1 \quad (4.11)$$

and hence

$$g = 0, \pm\frac{1}{2}, \pm 1, \dots, \quad (4.12)$$

which is precisely the Dirac condition. Equation (4.9) is called an 'equivariance' condition on the wave function Ψ : in going from the S_2 of (θ, ϕ) to the S_3 of (θ, ϕ, χ) we have enlarged

the configuration space over which our wave functions are to be defined, but an *arbitrary* dependence on the additional variable χ is not consistent with the required gauge invariance (dynamical independence) with respect to χ . Only Ψ 's satisfying (4.9) are allowed, with g satisfying (4.12). And, of course, our basic spinor

$$\begin{pmatrix} \cos \theta/2 & e^{i\chi} \\ \sin \theta/2 & e^{i(\phi+\chi)} \end{pmatrix} \quad (4.13)$$

does satisfy (4.12) with $g = -1/2$, the minimum non-trivial magnitude.

A general wave function $\Psi(s)$ can be expressed as a linear combination of the 'top' functions $\mathcal{D}_{m'm}^j(\theta, \phi, \chi)$, which carry irreducible representations of $SU(2)$. It seems obvious from the fact that θ , ϕ , and χ are angles that j should indeed be the angular momentum quantum number: for those who doubt, some further discussion is given in Appendix D. Then

$$\Psi(s) = \sum_{j,m',m} c_{m'm}^j \mathcal{D}_{m'm}^j(\theta, \phi, \chi). \quad (4.14)$$

The constraint (4.9) must now be imposed. If we multiply s from the right by $\exp(i\sigma_3\alpha)$, the \mathcal{D} 's get changed by

$$\mathcal{D}_{m'm}^j(s \exp i\sigma_3\alpha) = e^{2\alpha m i} \mathcal{D}_{m'm}^j(s), \quad (4.15)$$

since m is the eigenvalue of $\sigma_3/2$. Thus from (4.15) and (4.9),

$$m = g = 0, \pm \frac{1}{2}, \pm 1, \dots \quad (4.16)$$

and the possibility of 1/2-odd integral spin has emerged (since j is 1/2-odd integral if $2m$ is odd and integral if $2m$ is even). In fact, as stated in Section 1, the system has 1/2-odd angular momentum if the monopole strength g has the value $(n+1/2)$, for integer n .

5. The Skyrmiion case

Our basic analogy is as follows:

fast d.f.'s : fermion Fock states	~ spin states $ \uparrow\rangle, \downarrow\rangle$	}
fermion vacuum $ 0\rangle$	~ spin state $ \uparrow\rangle$	
slow d.f.'s: Goldstone boson fields $\phi_a \sim$ angular variables \mathbf{n}		}
$ 0, \phi_a\rangle \sim \uparrow \mathbf{n}\rangle$		

(5.1)

Just as a monopole structure appeared in the Berry phase associated with $|\uparrow, \mathbf{n}\rangle$ for slowly varying \mathbf{n} , so the W-Z term in the bosonic action is interpreted as a kind of Berry phase for $|0, \phi_a\rangle$.

We begin by introducing the commonly-used notation for the ϕ fields. In the case of $SU(2)_f$, we would have four ϕ 's, written as $\phi = (\sigma, \pi)$, where $\sigma^2 + \pi^2 = f^2$, quite analo-

224

gously to $n^2 = 1$. However, this does not generalize to the required $SU(3)_f$ case. Instead, we first rewrite ϕ as

$$\phi = fU, \quad (5.2)$$

where

$$U = \exp(i\tau \cdot \pi/f) \quad (5.3)$$

is a unitary 2×2 matrix. This amounts to a reparametrization of the original σ, π in the expression $\phi = (\sigma, \pi)$. In $SU(3)_f$, (5.3) is generalized to (cf. (1.22))

$$U = \exp(i\lambda \cdot \pi/f), \quad (5.4)$$

where π is understood now to be an 8-component ‘angle-type’ field. The analogue of (4.1) is then

$$\int \mathcal{D}U \exp\{i \int \mathcal{L}_0(U)dt\} \exp iS_{W-Z}(U), \quad (5.5)$$

where \mathcal{L}_0 is all the rest of the Lagrangian for the U fields, apart from the W–Z term; for example,

$$\mathcal{L}_0 = \frac{1}{4} \text{tr}(\partial_\mu U^\dagger \partial^\mu U) + \dots, \quad (5.6)$$

where the dots represent other terms which are necessary to stabilize the soliton, for instance. Finally, the expression for the W–Z action is [8]

$$\exp iS_{W-Z} = \exp \frac{-iN_c}{240\pi^2} \int \epsilon^{ijklm} \text{tr}(U^\dagger \partial_i U U^\dagger \partial_j U U^\dagger \partial_k U U^\dagger \partial_l U U^\dagger \partial_m U) d^5x. \quad (5.7)$$

The integral in (5.7) is over a 5-dimensional ‘disc’ whose boundary is 4-dimensional Minkowskian space-time. This disc is the 5-dimensional analogue of the 2-dimensional surfaces considered in (1.19)–(1.21), and N_c is the analogue [8] of the monopole N in (1.21).

Now, we seem a long way from anything like the Balachandran monopole Lagrangian (3.21). However, we can actually make the connection quite explicit, as follows [16, 17]. Instead of treating the full quantum-mechanical problem (5.5), in which the whole of the U matrix is treated as a quantum field variable, we perform only a ‘semi-classical’ quantization. In such an approach, one starts from a solution $U_c(r)$ of the static classical field equations in the $SU(2)_f$ case, which is of standard Skyrmion type (cf. (5.3) with $\pi = f\hat{r}\theta(r)$):

$$U_c(r) = \cos \theta(r) + i\tau \cdot \hat{r} \sin \theta(r). \quad (5.8)$$

It is clear that this solution is not rotationally invariant, nor is it invariant under isospin rotations. In fact, there are infinitely many such solutions, related to one another by spatial or isospin rotations, all of which are degenerate in energy since the original Lagrangian is invariant under space or isospin rotations. Actually these two kinds of rotation are effectively equivalent for (5.8), since

$$s\tau_i s^{-1} = \tau_j R_{ji}(s) \quad (5.9)$$

for $s \in \text{SU}(2)$. The coordinates which distinguish these degenerate classical configurations are the parameters of the matrix s . The semi-classical quantization procedure consists in promoting these d.f.'s into quantum variables $s(t)$. Thus we write

$$U(\mathbf{r}, t) = s(t)U_c(\mathbf{r})s^{-1}(t). \quad (5.10)$$

Classical quantities will now have a subscript c , and quantum d.f.'s will be distinguished by having no subscript c , instead of by having a '^'. The $s(t)$ will behave just like the s of Sections 3 and 4.

We must now extend (5.8) and hence (5.10) to the $\text{SU}(3)_f$ case, or else we get no $W - Z$ term at all [8] (see further Section 6). This means that we have to 'embed' (5.8) inside an $\text{SU}(3)$ matrix. The obvious way to do this would seem to be

$$U_c \rightarrow \begin{pmatrix} \cos \theta(r) + i\boldsymbol{\tau} \cdot \hat{\mathbf{r}} \sin \theta(r) & 0 \\ 0 & 1 \end{pmatrix} \equiv \tilde{U}_c(\mathbf{r}) \quad (5.11)$$

(alternative embeddings, which have different physical consequences, are discussed in Refs [12], [16] and [17]). So now,

$$U(\mathbf{r}, t) = s(t)\tilde{U}_c(\mathbf{r})s^{-1}(t), \quad s \in \text{SU}(3)_f. \quad (5.12)$$

We observe at once that $U(\mathbf{r}, t)$ is invariant under

$$s \rightarrow se^{iY\alpha(t)}, \quad (5.13)$$

where

$$Y = \frac{1}{3} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -2 \end{pmatrix} \quad (5.14)$$

(the normalization is, of course, chosen for convenience). Thus, the configuration space for the s d.f.'s is not $\text{SU}(3)$ but rather $\text{SU}(3)/\text{U}(1)_Y$; this is exactly analogous to our monopole example, where the required configuration space was $\text{SU}(2)/\text{U}(1)$. In fact, the analogy is very close indeed, for when (5.12) is inserted into (5.7) one finds — after some calculation — that the term involving s (i.e. the piece involving the quantum d.f.'s, in this approximation) is just [16, 17]

$$\mathcal{L}_{W-Z} = -\frac{1}{2} N_c B(U_c) \text{tr}(Y s^{-1} \dot{s}), \quad (5.15)$$

where $B(U_c)$ is the winding number (= baryon number) of the classical configuration U_c . Equation (5.15) should be compared with (3.21).

We see, from this comparison, that indeed the $W - Z$ term is acting so as to produce, in this semi-classical quantization, exactly a 'monopole in $\text{SU}(3)$ space'. The procedure of Section 4 can be transcribed easily to $\text{SU}(3)$. The gauge invariance analogous to (4.2) is the invariance of (5.13), under which, however, \mathcal{L}_{W-Z} changes according to

$$\mathcal{L}_{W-Z} \rightarrow \mathcal{L}_{W-Z} + \frac{1}{3} N_c B \dot{\alpha}. \quad (5.16)$$

In the quantized theory, s and Y are operators, and from Noether's theorem (corresponding to (4.7)) we have

$$\hat{Y}\Psi = \frac{1}{3}N_c B\Psi. \quad (5.17)$$

What is the analogue of the quantization constraint (4.12)? For this we note [17] that if we replace s by sh in (5.12), with $h \in \text{SU}(2)$, this is equivalent to rotating U_c by some spatial rotation parametrized by h (cf. 5.9)). In particular, consider a rotation by 2π about the 3rd axis. This corresponds to

$$h = \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} \quad (5.18)$$

and thus to the replacement

$$s \rightarrow s \begin{pmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = se^{3\pi i Y}. \quad (5.19)$$

According to (5.17), the allowed Ψ 's must then pick up a phase factor

$$e^{i\pi N_c B}, \quad (5.20)$$

and hence *the allowed states for $B = 1$ are fermions if N_c is odd, and bosons if N_c is even!* [9].

The wave functionals Ψ are the $\text{SU}(3)$ generalization of the $\text{SU}(2)$ rotation functions $\mathcal{D}_{m'm}^j$ [$s \in \text{SU}(2)$] — namely

$$\mathcal{D}_{I,I_3,Y;I',I'_3,Y'}^{p,q} (s \in \text{SU}(3)), \quad (5.21)$$

where p and q label the irreducible representation of $\text{SU}(3)$. In (5.21) the left-hand group of ‘magnetic quantum numbers’ refers to transformation properties under *left* multiplication of s by a matrix in $\text{SU}(3)$, and hence (cf. (5.12)) to a flavour rotation of U ; the right-hand indices refer to right multiplication. But we have already seen that the $\text{SU}(2)$ part — in the sense of (5.11) — of any ‘right multiplication’ matrix corresponds to a spatial rotation. Hence I' and I'_3 are actually the real spin and its third component. Now for $B = 1$ and $N_c = 3$ we need the eigenvalue $Y' = 1$ from (5.17): The lowest dimensionality $\text{SU}(3)$ representations with $Y' = 1$ are the **8** and **10** (Fig. 3). In the former, the states with $Y' = 1$ have $I' = 1/2$, and hence spin 1/2, while in the latter they have spin 3/2. The left-hand indices give just the flavour quantum numbers corresponding to these $\text{SU}(3)$ representations: thus we have an **8** of spin 1/2 and a **10** of spin 3/2.

Further details of Skyrme quantization are given in Guadagnini [11] and Rabino-vici et al. [12]: our concern here has been to place the ‘monopole’ form (5.15) of $\mathcal{L}_{w,z}$ in the context of an adiabatic decoupling problem. From this point of view, the peculiar phase behaviour leading to ‘fermion-ness’ in the ϕ sector has arisen as a result of non-trivial structure left behind when the fermion vacuum is decoupled adiabatically from

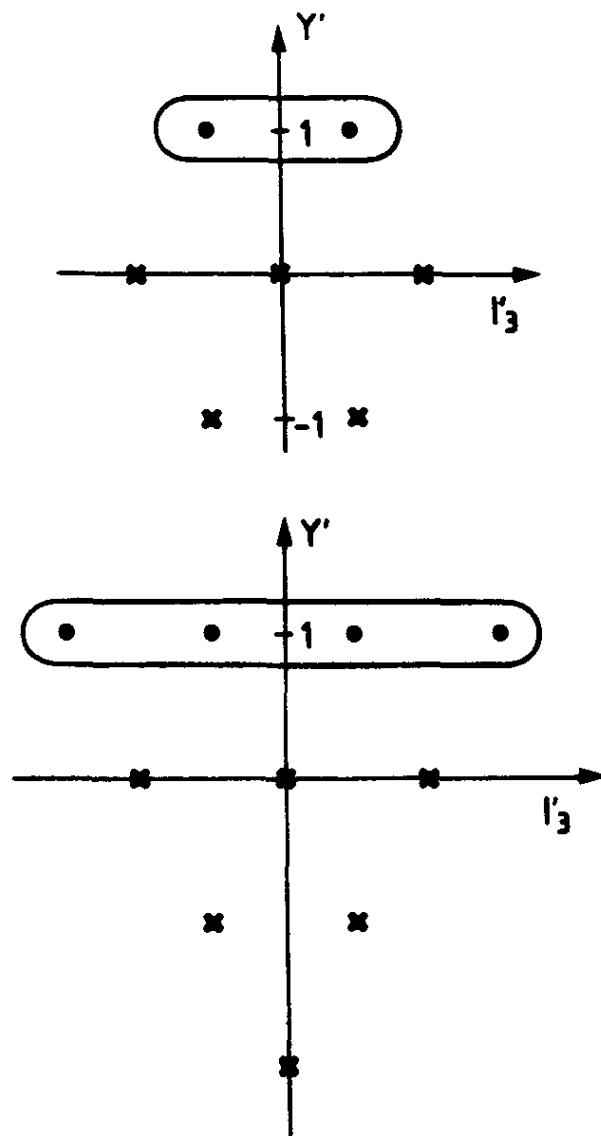


Fig. 3. The **8** and **10** representations of $SU(3)$, showing the two allowed multiplets with $Y' = 1$

the ϕ 's. If we use only the ϕ d.f.'s, and integrate the fermions away, we must include a $W - Z$ term which embodies this structure. The ultimate reason that this structure has a 'monopole' form is to be found in the topological approach to anomalies [27, 28].

We may also remark that a similar mechanism holds for Skyrmions in $2+1$ dimensions [33]. Here the Wess-Zumino term is replaced by the Hopf term [34-37], which can also be interpreted as arising from integrating out fermions [35]. When the Skyrmion is quantized semi-classically, the angular momentum has the value (*integer* + $\theta/2\pi$), where θ is the coefficient of the Hopf term [38]. The effective Lagrangian in this approximation is exactly analogous to that describing a charged particle moving in two dimensions in the field of a magnetic vortex lying perpendicular to the plane of motion. In that case, the angular momentum can have a value neither integral nor half-odd integral, with consequential 'fractional statistics' [39-41]. Thus just as, in $3+1$, the $W - Z$ term acts as a monopole in field space, so in $2+1$ the Hopf term acts as a vortex in field space.

6. Postscript: the case of only two flavours

The above discussion has been predicated upon the existence of the $W-Z$ term — whose presence determines the quantization of the Skyrmi n (fermion if N_c odd, boson if N_c even). But if there are only two flavours, there is no $W-Z$ term: when (5.3) is substituted into (5.7) the $SU(2)$ trace vanishes¹. Yet there are topological solitons since $\pi_3(SU(2)) = \mathbb{Z}$. What determines their quantization?

The answer is that the $B = 1$ soliton can be quantized *either* as a boson *or* as a fermion: there is no restriction involving N_c , and one has to choose the fermionic option by hand [42]. The way in which fermionic quantization is *possible* (but not required) was discussed by Finkelstein [43], Finkelstein and Rubinstein [44], and Williams [45]. One way of putting it is as follows [46]. Since $\pi_4(SU(2)) = \mathbb{Z}_2$, time-dependent soliton fields U fall into two distinct homotopy classes of maps from (compactified) space-time to $SU(2)$. Functional integrals over the U 's can therefore be separated into two topologically disjoint sectors (analogous to θ -vacua in QCD), corresponding to those U 's which can be continuously deformed to the identity, and those which cannot. The contribution from these two sectors to the functional propagator can have a relative + sign or a relative - sign: in the former case the propagator contains all integral spins (bosonic), in the latter half-integral ones (fermionic).

This situation is mathematically the same as that of the spherical top [3], since $\pi_4(O(3)) = \mathbb{Z}_2$ also, and the same boson/fermion option therefore exists. In this case one can say, alternatively, that since $O(3)$ is doubly-connected, wave functions on $O(3)$ need not be single-valued. One can define single-valued wave functions by passing to the universal covering space $SU(2)$, but then one has to project back to $O(3)$ via $SU(2) \rightarrow O(3) \simeq SU(2)/\mathbb{Z}_2$, on which a double-valuedness can appear.

For $N_f > 2$, $\pi_4(SU(N_f)) = 0$ and so this boson/fermion *option* is removed. But then, since $\pi_5(SU(N_f)) = \mathbb{Z}$ we have the $W-Z$ addition to the Lagrangian, and the N_c -related quantization is determined.

I am grateful to Jo Zuk for many very helpful discussions; and to Stephen Wilkinson for patient instruction in some of the relevant mathematics, and for carefully reading the manuscript. It is a pleasure to take this opportunity of thanking Drs M. Praszałowicz and W. Słomiński for organising such a stimulating and enjoyable School, and for their warm hospitality.

APPENDIX A

Monopole strength and $(U)1$ winding number

We have seen that the monopole strength g in

$$\mathcal{L}_{\text{eff}}(n) = gi \operatorname{tr}(\sigma_3 s^{-1} \dot{s}) \quad (\text{A.1})$$

¹ Alternatively [8], in $SU(2)$ G -parity invariance forbids amplitudes with an odd number of pions, while (5.7) would, if it were non-vanishing, allow them; in $SU(3)$, (5.7) allows $K\bar{K} \rightarrow 3\pi$, which is not forbidden by G -parity.

is restricted to the values $g = p/2$ where $p = 0, \pm 1, \pm 2, \dots$. In this Appendix we will show how p can be interpreted as a winding number associated with the $U(1)$ gauge transformation (4.2).

Consider a *sequence* of gauge transformations

$$s \rightarrow s \exp i\sigma_3 \alpha(t) \quad (\text{A.2})$$

parametrized by t , where

$$\alpha(t=0) = 0, \quad \alpha(t=T) = 2\pi, \quad (\text{A.3})$$

so that we have a closed loop in s -space,

$$s(t=0) = s(t=T). \quad (\text{A.4})$$

Then as we move through this sequence of t -values, the parameter α of the $U(1)$ gauge group goes once round its circle (Fig. A1a).

Corresponding to the gauge transformation (A.2), we have the transformation

$$z = \begin{pmatrix} \cos \theta/2 e^{ix} \\ \sin \theta/2 e^{i(\phi+x)} \end{pmatrix} \rightarrow e^{i\alpha(t)} z \quad (\text{A.5})$$

of the basic \uparrow spinor (cf. (3.6) and (3.12)). Thus $\alpha(t)$ is just a variable phase for the associated spinor, and as we go round the sequence of gauge transformations in Fig. A1a, this phase swings round precisely once (Fig. A1b). Meanwhile, what is happening to $\mathcal{L}_{\text{eff}}(n)$? This becomes

$$\mathcal{L}_{\text{eff}}(n) \rightarrow \mathcal{L}_{\text{eff}}(n) - p\dot{\alpha}, \quad (\text{A.6})$$

where $p = 2g$. The associated effective *action* therefore changes by

$$\exp -i \int_0^T p \dot{\alpha} dt = \exp (-2\pi i p), \quad (\text{A.7})$$

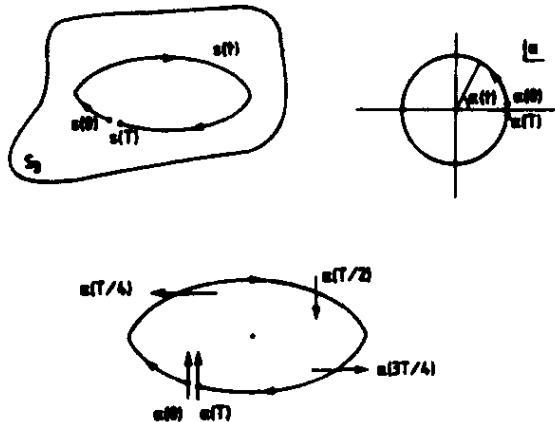


Fig. A1. Sequence of gauge transformations corresponding to a closed loop in s -space, and the associated variation of the spinor phase

230

i.e. its phase swings round p times as we follow the circuit of Fig. A1a. We can therefore interpret p as a winding number which counts the number of rotations of the action phase as we circulate once in α -space (i.e. one circuit in $U(1)$ space).

APPENDIX B

More on the Hopf map

In Section 3 we gave two forms of the Hopf map, one in terms of z

$$\mathbf{n} = z^\dagger \boldsymbol{\sigma} z, \quad (\text{B.1})$$

and the other in terms of s

$$\boldsymbol{\sigma} \cdot \mathbf{n} = s \sigma_3 s^{-1}, \quad (\text{B.2})$$

where

$$z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}, \quad |z_1|^2 + |z_2|^2 = 1 \quad (\text{B.3})$$

and

$$s = \begin{pmatrix} z_1 & -z_2^* \\ z_2 & z_1^* \end{pmatrix}. \quad (\text{B.4})$$

We make the connection between (B.1) and (B.2) as follows. Let us write $z_1 = x_1 + ix_2$, $z_2 = x_3 + ix_4$. Then from (B.1)

$$\mathbf{n}_1 = (x_1 - ix_2 \quad x_3 - ix_4) \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} x_1 + ix_2 \\ x_3 + ix_4 \end{pmatrix} = 2(x_1 x_3 + x_2 x_4) \quad (\text{B.5})$$

and

$$\mathbf{n}_2 = 2(x_1 x_4 - x_2 x_3) \quad (\text{B.6})$$

$$\mathbf{n}_3 = x_1^2 + x_2^2 - x_3^2 - x_4^2, \quad (\text{B.7})$$

while

$$x_1^2 + x_2^2 + x_3^2 + x_4^2 = 1. \quad (\text{B.8})$$

On the other hand, for our s of (3.12), connected to the parameters θ, ϕ of \mathbf{n} , we had

$$z_1 = \cos \theta/2 e^{i\chi}; x_1 = \cos \theta/2 \cos \chi, \quad x_2 = \cos \theta/2 \sin \chi \quad (\text{B.9})$$

$$z_2 = \sin \theta/2 e^{i(\phi+\chi)}; x_3 = \sin \theta/2 \cos(\phi + \chi), \quad x_4 = \sin \theta/2 \sin(\phi + \chi), \quad (\text{B.10})$$

whence from (B.5)–(B.7)

$$\left. \begin{aligned} \mathbf{n}_1 &= \sin \theta \cos \phi \\ \mathbf{n}_2 &= \sin \theta \sin \phi \\ \mathbf{n}_3 &= \cos \theta \end{aligned} \right\} \quad (\text{B.11})$$

as required.

The matrix s has a simple geometrical interpretation. Consider first the case of

$$s_+(\theta, \phi) = \begin{pmatrix} \cos \theta/2 & -\sin \theta/2 e^{-i\phi} \\ \sin \theta/2 e^{i\phi} & \cos \theta/2 \end{pmatrix}. \quad (\text{B.12})$$

Let \hat{u} be the unit vector

$$\hat{u} = (-\sin \phi, \cos \phi, 0) \quad (\text{B.13})$$

and consider

$$\exp(-i\sigma \cdot \hat{u}\theta/2) = \cos \theta/2 - i \sin \theta/2 \sigma \cdot \hat{u} \quad (\text{B.14})$$

$$= s_+(\theta, \phi). \quad (\text{B.15})$$

This is a rotation of θ about \hat{u} , which rotates \hat{z} into \hat{n} (Fig. B1). So Eq. (B.2), which is equivalent to

$$s_+^{-1} \sigma \cdot n s_+ = \sigma_3 \quad (\text{B.16})$$

in this case, means simply that n has been rotated to be along the 3rd axis. The remaining factor $\exp(i\sigma_3\chi)$ in (3.29) is then just a rotation about the 3 axis (the 'body' axis).

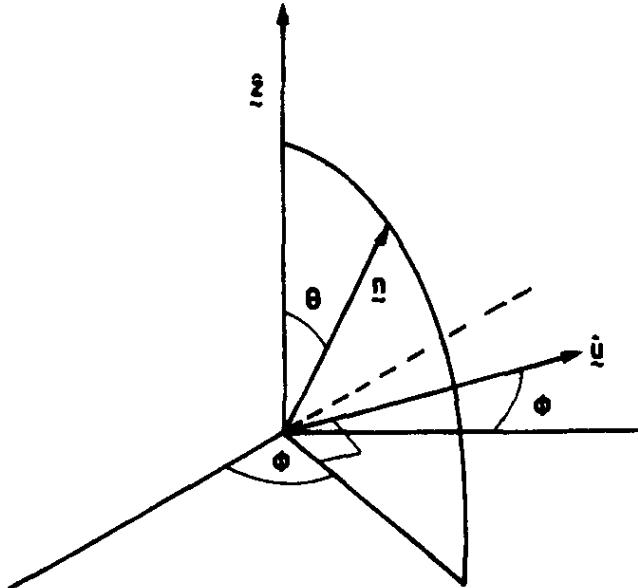


Fig. B1. A rotation of θ about \hat{u} rotates \hat{z} into \hat{n}

There is yet one more way of writing the connection between z (or s) and n which the Hopf map enforces. It is

$$\sigma \cdot n = 2zz^\dagger - 1. \quad (\text{B.17})$$

This is easy to verify: using (B.1) for the components of n in terms of z_1, z_2 , we find

$$\sigma \cdot n = \begin{pmatrix} |z_1|^2 - |z_2|^2 & 2z_2^* z_1 \\ 2z_1^* z_2 & |z_2|^2 - |z_1|^2 \end{pmatrix} = 2zz^\dagger - 1, \quad (\text{B.18})$$

with the help of $|z_1|^2 + |z_2|^2 = 1$.

APPENDIX C

An alternative quantization for monopoles [12]

Let us consider the monopole action in (4.1),

$$S = \int [\frac{1}{2} \dot{\mathbf{n}}^2 + \frac{1}{2} i p \operatorname{tr} (\sigma_3 s^{-1} \dot{s})] dt \quad (\text{C.1})$$

$$= \int [\frac{1}{2} \dot{\mathbf{n}}^2 + i p z^\dagger \dot{z}] dt, \quad (\text{C.2})$$

where the second step follows from (3.21), and g has been replaced by $p/2$ (Appendix A). We have not so far considered explicitly the *first* ('kinetic energy') term of (C.2) — let us attend to it now.

We have

$$\frac{1}{2} \dot{\mathbf{n}}^2 = \frac{1}{4} \operatorname{tr} (\boldsymbol{\sigma} \cdot \dot{\mathbf{n}})^2 \quad (\text{C.3})$$

$$= 2[(\dot{z}^\dagger \dot{z}) + (z^\dagger \dot{z})^2], \quad (\text{C.4})$$

using (B.17). Let us write

$$a = iz^\dagger \dot{z}; \quad (\text{C.5})$$

then

$$\frac{1}{2} \dot{\mathbf{n}}^2 = 2\dot{z}^\dagger \dot{z} - 2a^2. \quad (\text{C.6})$$

But also

$$\operatorname{tr} (\dot{s}^\dagger \dot{s}) = 2\dot{z}^\dagger \dot{z}, \quad (\text{C.7})$$

by direct verification. Hence finally we can write this part of the action as

$$\int [\operatorname{tr} (\dot{s}^\dagger \dot{s}) - 2a^2] dt. \quad (\text{C.8})$$

Consider now the behaviour of (C.8) under the U(1) gauge transformation:

$$z \rightarrow e^{i\alpha(t)} z, \quad (\text{C.9})$$

which also corresponds to

$$s \rightarrow se^{i\alpha(t)\sigma_3}. \quad (\text{C.10})$$

Under (C.9),

$$a \rightarrow a - \dot{\alpha}, \quad (\text{C.11})$$

so that

$$-2a^2 \rightarrow -2a^2 + 4a\dot{\alpha} - 2\dot{\alpha}^2. \quad (\text{C.12})$$

On the other hand, under (C.10) one finds easily

$$\operatorname{tr} (\dot{s}^\dagger \dot{s}) \rightarrow \operatorname{tr} (\dot{s}^\dagger \dot{s}) - 4a\dot{\alpha} + 2\dot{\alpha}^2. \quad (\text{C.13})$$

Thus (C.9) and (C.10) are together an invariance of (C.8), as we stated in Section 4.

We can bring this invariance out by rewriting (C.8) as

$$\int \operatorname{tr} \{ [(\vec{\partial}_t - ia\sigma_3)s^\dagger] [s(\vec{\partial}_t + ia\sigma_3)] \} \quad (\text{C.14})$$

as can be simply checked, recalling that $a = \frac{1}{2} \text{tr}(\sigma_3 s^{-1} \dot{s})$ also. Expression (C.14) is manifestly invariant under the combined transformations (C.9) and (C.10). Thus we are interested in the generating functional

$$Z = \int \mathcal{D}s \exp i[\oint (\text{tr} \{[(\vec{\partial}_t - ia\sigma_3)s^\dagger] [s(\vec{\partial}_t + ia\sigma_3)]\} + pa)dt], \quad (\text{C.15})$$

where the action is evaluated over closed loops in s -space. This can be rewritten with the aid of an auxiliary field $A(t)$ [12] as

$$Z \sim \int \mathcal{D}A \mathcal{D}s \exp i[\oint (\text{tr} \{[(\vec{\partial}_t - iA\sigma_3)s^\dagger] [s(\vec{\partial}_t + iA\sigma_3)]\} + pA + \frac{1}{2}(p/2)^2)dt], \quad (\text{C.16})$$

where a constant

$$\int \mathcal{D}A \exp \{2[A + (\frac{1}{4}p - a)]^2\}$$

has been ignored.

In (C.16) A acts as an independent gauge field, which changes by $A \rightarrow A - \dot{\alpha}$ under the gauge transformation (C.10), so that (C.16) is gauge invariant. We can work in the specific gauge $A = 0$, and require that the equation of motion obtained from the variation with respect to A (i.e. Gauss's law for this case) be realized as a constraint on the physical states. In this gauge the Lagrangian of (C.16) is just

$$\mathcal{L}(A = 0) = \text{tr}(\dot{s}^\dagger \dot{s}) + \frac{1}{2}(p/2)^2 \quad (\text{C.17})$$

and Gauss's law is

$$\text{tr}(-i\sigma_3 s^\dagger \dot{s} + i\dot{s}^\dagger s\sigma_3) = -p. \quad (\text{C.18})$$

Equation (C.18) is the equivalent of (4.7), since the l.h.s. can be identified with the generator of right transformations (C.10), as we discuss further in Appendix D. Indeed, as we also show there, the term $\text{tr}(\dot{s}^\dagger \dot{s})$ is precisely $\frac{1}{2}\mathbf{J}^2$, the square of the angular momentum operator (the motion in r being ignored, only the angles varying). The Hamiltonian in this gauge is therefore

$$H(A = 0) = \frac{1}{2}\mathbf{J}^2 - \frac{1}{2}(p/2)^2 \quad (\text{C.19})$$

and the eigenfunctions are again $\mathcal{D}_{m,m}^j(s)$ with m (which carries the right multiplications) restricted to the value $-p/2, p = 0, \pm 1, \pm 2, \dots$. The eigenvalues are $\frac{1}{2}(j(j+1) - (p/2)^2)$, the allowed j being $j = |p/2|, |p/2| + 1, \dots$.

APPENDIX D

Angular momentum

The connection between the s -formalism and the conventional 'spherical top' formalism can be made explicit by parametrizing s by the Euler angles α, β, γ according to

$$s = \begin{pmatrix} e^{i(\alpha+\gamma)/2} \cos \beta/2 & e^{i(\gamma-\alpha)/2} \sin \beta/2 \\ -e^{-i(\gamma-\alpha)/2} \sin \beta/2 & e^{-i(\alpha+\gamma)/2} \cos \beta/2 \end{pmatrix}. \quad (\text{D.1})$$

234

Straightforward calculation then yields

$$T \equiv \text{tr}(\dot{s}^\dagger s) = \frac{1}{2}\dot{\beta}^2 + \frac{1}{2}(\dot{\gamma} + \dot{\alpha} \cos \beta)^2 + \frac{1}{2}\dot{\alpha}^2 \sin^2 \beta, \quad (\text{D.2})$$

which may be compared with the expression for the spherical top kinetic energy given by Edmonds [47], p. 66. The momenta canonically conjugate to α , β , γ are then

$$p_\alpha = \frac{\partial T}{\partial \dot{\alpha}} = \dot{\alpha} + \dot{\gamma} \cos \beta, \quad \text{etc.,} \quad (\text{D.3})$$

and the passage to quantum theory is made by

$$p_\alpha \rightarrow -i \frac{\partial}{\partial \alpha}, \quad \text{etc.} \quad (\text{D.4})$$

We find

$$\begin{aligned} T = & -\frac{1}{2} \left\{ \frac{\partial^2}{\partial \beta^2} + \cot \beta \frac{\partial}{\partial \beta} + \cosec^2 \beta \frac{\partial^2}{\partial \gamma^2} \right. \\ & \left. + \frac{1}{\sin^2 \beta} \frac{\partial^2}{\partial \alpha^2} - \frac{2 \cos \beta}{\sin^2 \beta} \frac{\partial^2}{\partial \alpha \partial \gamma} \right\} \end{aligned} \quad (\text{D.5})$$

for the operator representing the (rotational) kinetic energy. This is, in fact, precisely the angular kinetic energy

$$T = \frac{1}{2} \mathbf{J}^2 \quad (\text{D.6})$$

following Edmonds [47]. The eigenfunctions of T are then $\mathcal{D}_{m'm}^j(\alpha, \beta, \gamma)$. Alternative parametrizations of s , such as (3.29), are, of course, also possible.

Finally, we note that

$$i \text{tr}(\dot{s}^\dagger s \sigma_3 - \sigma_3 s^\dagger \dot{s}) = 2(\dot{\alpha} + \dot{\gamma} \cos \beta) = 2p_\alpha. \quad (\text{D.7})$$

In the quantum theory, p_α is the generator of rotations about the 3-axis, which are represented in terms of s by right transformations

$$s(\alpha, \beta, \gamma) = s(\alpha = 0, \beta, \gamma) e^{i\sigma_3 \alpha/2}. \quad (\text{D.8})$$

Thus $2p_\alpha$ is the generator associated with the transformation (C.10), as claimed in Appendix C, and its eigenvalues should be integral — as indeed is required by the constraint (C.18).

REFERENCES

- [1] T. H. R. Skyrme, *Proc. R. Soc. London A* **260**, 127 (1961).
- [2] F. Bopp, R. Haag, *Z. Naturforsch.* **5a**, 644 (1950).
- [3] L. S. Schulman, *Phys. Rev.* **176**, 1558 (1968); see also L. S. Schulman, *J. Math. Phys.* **12**, 304 (1971).
- [4] P. A. M. Dirac, *Proc. R. Soc. London A* **133**, 60 (1931).
- [5] I. Tamm, *Z. Phys.* **71**, 141 (1931).

- [6] R. Jackiw, C. Rebbi, *Phys. Rev. Lett.* **36**, 1116 (1976); P. Hasenfratz, G. 't Hooft, *Phys. Rev. Lett.* **36**, 1119 (1976).
- [7] A. Goldhaber, *Phys. Rev. Lett.* **36**, 1122 (1976).
- [8] E. Witten, *Nucl. Phys.* **B223**, 422 (1983).
- [9] E. Witten, *Nucl. Phys.* **B223**, 433 (1983).
- [10] J. Wess, B. Zumino, *Phys. Lett.* **37B**, 95 (1971).
- [11] E. Guadagnini, *Nucl. Phys.* **B236**, 35 (1984).
- [12] E. Rabinovici, A. Schwimmer, S. Yankelowicz, *Nucl. Phys.* **B248**, 523 (1984).
- [13] A. P. Balachandran, G. Marmo, B.-S. Skagerstam, A. Stern, *Gauge symmetries and fibre bundles*, Springer-Verlag, Berlin etc. 1983.
- [14] A. P. Balachandran, G. Marmo, B.-S. Skagerstam, A. Stern, *Nucl. Phys.* **B162**, 385 (1980).
- [15] F. Zaccaria et al., *Phys. Rev.* **D27**, 2327 (1983).
- [16] A. P. Balachandran, F. Lizzi, V. G. J. Rodgers, A. Stern, *Nucl. Phys.* **B256**, 525 (1985).
- [17] A. P. Balachandran, TASI lectures at Yale University, 9 June-15 July 1985.
- [18] M. V. Berry, *Proc. R. Soc. London* **A392**, 45 (1984).
- [19] M. Stone, Illinois preprint ILL-(TH)-85-#55 (1985).
- [20] Y. S. Wu, A. Zee, *Nucl. Phys.* **B258**, 157 (1985).
- [21] Y. S. Wu, A. Zee, Santa Barbara preprint NSF-ITP-85-128 (1985).
- [22] I. J. R. Aitchison, C. M. Fraser, *Phys. Rev.* **D31**, 2605 (1985).
- [23] P. Simic, *Phys. Rev. Lett.* **55**, 40 (1985); and Rockefeller University preprint RU84/B/106 (1984).
- [24] H. Kuratsuji, S. Iida, *Prog. Theor. Phys.* **74**, 439 (1985).
- [25] A. Messiah, *Quantum mechanics*, North Holland, Amsterdam 1962, Vol. II, pp. 781-800.
- [26] F. Wilczek, A. Zee, *Phys. Rev. Lett.* **52**, 2111 (1984).
- [27] L. Alvarez-Gaumé, P. Ginsparg, *Nucl. Phys.* **B243**, 449 (1984).
- [28] P. Nelson, L. Alvarez-Gaumé, *Commun. Math. Phys.* **99**, 103 (1985).
- [29] T. T. Wu, C. N. Yang, *Phys. Rev.* **D12**, 3845 (1975).
- [30] T. T. Wu, C. N. Yang, *Nucl. Phys.* **B107**, 365 (1976).
- [31] L. H. Ryder, *J. Phys. A: Math. Gen.* **13**, 437 (1980).
- [32] M. Minami, *Prog. Theor. Phys.* **62**, 1128 (1979).
- [33] A. A. Belavin, A. M. Polyakov, *JETP Lett.* **22**, 245 (1975).
- [34] F. Wilczek, A. Zee, *Phys. Rev. Lett.* **51**, 2250 (1983).
- [35] A. M. Din, W. J. Zakrzewski, *Phys. Lett.* **146B**, 341 (1984).
- [36] Y. S. Wu, A. Zee, *Phys. Lett.* **147B**, 325 (1984).
- [37] T. Jaroszewicz, *Phys. Lett.* **146B**, 337 (1984); Harvard preprint HUTP/BO1 (1985).
- [38] M. J. Bowick, D. Karabali, L. C. R. Wijewardhana, Yale University preprint YTP 85-20 (1985).
- [39] F. Wilczek, *Phys. Rev. Lett.* **48**, 1144 (1982).
- [40] F. Wilczek, *Phys. Rev. Lett.* **49**, 957 (1982).
- [41] R. Jackiw, A. N. Redlich, *Phys. Rev. Lett.* **50**, 555 (1983).
- [42] G. S. Adkins, C. R. Nappi, E. Witten, *Nucl. Phys.* **B228**, 552 (1983).
- [43] D. Finkelstein, *J. Math. Phys.* **7**, 1218 (1966).
- [44] D. Finkelstein, J. Rubinstein, *J. Math. Phys.* **9**, 1762 (1968).
- [45] J. W. Williams, *J. Math. Phys.* **11**, 2611 (1970).
- [46] J. S. Dowker, *J. Phys. A*, Ser. 2, No. **5**, 936 (1972); see also J. S. Dowker, *Lett. Nuovo Cimento* **4**, 301 (1972), and L. S. Schulman, *J. Math. Phys.* **12**, 304 (1971).
- [47] A. R. Edmonds, *Angular momentum in quantum mechanics*, Princeton University Press, Princeton N. J. 1957.

Hamiltonian Interpretation of Anomalies

Philip Nelson^{1*} and Luis Alvarez-Gaumé²

¹ Institute for Theoretical Physics, University of California, Santa Barbara, CA 93106, USA

² Lyman Laboratory of Physics, Harvard University, Cambridge, MA 02138, USA

Abstract. A family of quantum systems parametrized by the points of a compact space can realize its classical symmetries via a new kind of nontrivial ray representation. We show that this phenomenon in fact occurs for the quantum mechanics of fermions in the presence of background gauge fields, and is responsible for both the nonabelian anomaly and Witten's SU(2) anomaly. This provides a hamiltonian interpretation of anomalies: in the affected theories Gauss' law cannot be implemented. The analysis clearly shows why there are no further obstructions corresponding to higher spheres in configuration space, in agreement with a recent result of Atiyah and Singer.

1. Introduction

We say we have an “anomaly” when a symmetry of a classical field theory is not reflected at all in those of the corresponding quantum theory, or more precisely when the full set of classical symmetries cannot be preserved in any of the many possible quantization schemes. When the symmetry in question is an ordinary one such as scale or chiral invariance, we have a straightforward interpretation for the effects of the anomaly in terms of states in Hilbert space: the symmetry in question is absent from the full theory. Coupling constants run; tunneling events do not conserve axial charge. These results are surprising, but not fatal to the theory.

The case of gauged symmetries is very different. Gauge symmetries are properly to be thought of as not being symmetries at all, but rather redundancies in our description of the system [1]. The true configuration space of a (3+1)-dimensional gauge theory is the quotient $\mathcal{C}^3 = \mathcal{A}^3 / \mathcal{G}^3$ of gauge potentials in $A_0 = 0$ gauge modulo three-dimensional gauge transformations¹. When gauge degrees of freedom become anomalous, we find that they are not redundant after all.

* Harvard Society of Fellows. Permanent address: Lyman Laboratory of Physics, Harvard University, Cambridge, MA 02138, USA

¹ We will sometimes omit the superscript 3

Recently it has become clear that gauge theories with fermion display three different kinds of anomalies, all related to the global topology of the four-dimensional configuration space \mathcal{C}^4 by the family index of the Dirac operator \not{D}^4 . These are the axial U(1) anomaly [the “ $\pi_0(\mathcal{G}^3)$ anomaly”], Witten’s SU(2) anomaly [2] [from $\pi_1(\mathcal{G}^3)$], and the nonabelian gauge anomaly [3] [from $\pi_2(\mathcal{G}^3)$]. The diversity of the manifestations of these anomalies seems to belie their common origin, however. In the first case we find particle production in the presence of instanton fields [4], breaking of a global symmetry, and no problem with gauge invariance. In the second we find no problem with chiral charge, but instead a nonperturbative failure of *gauge* symmetry, while in the latter the same thing occurs even perturbatively.

What is going on? In the following sections we will attempt to give a hamiltonian picture of the gauge anomalies as simple as the axial anomaly’s particle-production interpretation. Essentially the answer will be that in anomalous theories we cannot formulate any Gauss law to constrain the physical states. Along the way we will try to make the above differences a bit less mystifying than they seem in the lagrangian picture. They will all turn out merely to reflect a simple fact about codimension: removing a point from a manifold can sever it into disconnected pieces only if its dimension equals one.

The aim of this paper is expository. We will not find any previously unknown anomalies, but instead will give an approach to understanding them which we have found illuminating. Our point of departure was a remark in [2] which we have generalized to embrace the anomaly of [3] as well². In Sect. 2 we set up our framework and establish our criterion for a global anomaly to exist. In Sects. 3 and 4 we verify the criterion for the cases of [2, 3] respectively, making use of known results from the lagrangian approach. In Sect. 5 we conclude with remarks.

2. Setting Up

It may seem difficult to arrive at a physical interpretation of a problem which renders a gauge theory nonsensical. We know, however, that anomalies do not themselves originate in the gauge sector. We can therefore attempt to quantize a given theory in two steps, starting with the matter fields; at the intermediate point we will have a *family* of quantum systems parametrized by the space of classical background gauge field configurations \mathcal{A}^3 . Furthermore, the whole collection should realize the classical gauge symmetry *via* unitary operators. The situation is not quite like the usual case of symmetry in quantum mechanics [6], however, since the transformations in question act both on Hilbert space \mathcal{H} and on background configuration space \mathcal{A} . They are indeed *bundle maps* of a *family* of Hilbert spaces, $\mathcal{H} \xrightarrow{\pi} \mathcal{A}$. A simple example of such a situation is an ordinary quantum mechanics problem with a Schrödinger particle interacting with a classical rotor degree of freedom $\bar{\phi}$: for fixed position of the rotor the system has no rotational symmetry, but the full family of theories does have an invariance expressed as a set of isometries, $U_\alpha : \mathcal{H}_{\bar{\phi}} \rightarrow \mathcal{H}_{\bar{\phi} + \alpha}$.

² We were also influenced by the work of Rajeev [5]

The notion of families of quantum systems has recently appeared in several papers [7–9]. The phenomenon of “quantum holonomy” discussed in these papers will be crucial to our analysis.

When an ordinary quantum system realizes its classical symmetries, however, it need not do so in the obvious way, by a unitary action of the symmetry group G on \mathcal{H} . Instead, Wigner showed [6] that in general we can demand only that \mathcal{H} furnish a projective, or ray, representation of G . When G is a multiply-connected topological group, \mathcal{H} will thus in general have irreducible sectors transforming under \tilde{G} , the universal cover of G . This is, of course, the situation with the rotation group, where \mathcal{H} has a sector of odd fermion number transforming as a “double-valued representation of $O(3)$,” i.e., as a true representation of spin (3).

The same sort of thing can occur in parametrized families of quantum systems. As a simple example, let us return to the case of the Schrödinger particle and rotor. Constraining the particle to lie on a circle, we have the hamiltonian

$$H = -(\partial_\phi)^2 + b[\delta(\phi - \frac{1}{2}\bar{\phi}) + \delta(\phi - \frac{1}{2}\bar{\phi} + \pi)], \quad (1)$$

which is continuous for $\bar{\phi} \in S^1$. For each $\bar{\phi}$ half of the energy eigenstates of this system are odd under the translation $\phi \rightarrow \phi + \pi$. Now let $\bar{\phi}$ vary, and for each value choose a real energy eigenfunction $\psi_{\bar{\phi}}$ with fixed eigenvalue ϵ . For the odd states it will be impossible to choose $\psi_{\bar{\phi}}$ smoothly; as $\bar{\phi}$ completes a full circuit ψ goes over to its negative. In other words, the odd energy eigenspaces each form twisted line bundles over the parameter space S^1 .

Let us attempt to find a unitary action of the symmetry group $U(1)$ on a given odd energy eigenspace \mathcal{H}_n of \mathcal{H} . Clearly U_g must map \mathcal{H}_n to itself, but at each point a decision must be made: there is no canonical choice of sign. This raises the possibility that *no* smooth choice may exist. Indeed, any ordinary unitary action of the symmetry group $U(1)$ must take any given ψ at $\bar{\phi}=0$ and give a nonzero section of \mathcal{H} . Since no such section exists, this quantum system cannot realize its $U(1)$ symmetry via an ordinary unitary action.³ More formally, if the Hilbert bundle \mathcal{H} admits an action of $G = U(1)$ which projects to the usual action of $U(1)$ on the parameter space S^1 , we say it is a “ G -bundle” [10]. In this case \mathcal{H} reduces to a new bundle $\bar{\mathcal{H}}$ defined on the quotient $S^1/U(1) = \text{point}$, and so is trivial. That is, any nontrivial bundle on the base (in our case an energy eigenspace) is not a G -bundle.

If the parameter space consists of many G -orbits it is sufficient to show that any one is nontrivial in order to rule out an ordinary G -action. In any case the key feature which makes possible the unremovable minus sign in the group action is the fact that the orbits are copies of G , which is not simply-connected.

Suppose now that we wish to quantize the rotor degree of freedom as well. The wavefunctions of the complete system can then be taken as complex functions of both ϕ , the particle position, and $\bar{\phi}$, the rotor position. Alternatively, however, they can be taken as functions from S^1 into the *space of functions* of ϕ , that is, as *sections* of the Hilbert bundle \mathcal{H} . We will call the complete Hilbert space

³ The reader may well object that we have simply chosen a foolish normalization for the $U(1)$ generator. Indeed the model has another classical $\overline{U(1)}$ symmetry which is realized in the usual way. We will return to this point.

$\dot{\mathcal{H}} \equiv \Gamma(\mathcal{H})$, the space of sections. It has a subspace spanned by the even eigenfunctions, and on this subspace we can define the unitary operator \mathcal{U}_α by $\mathcal{U}_\alpha \psi = \psi'$, where $\psi'_\varphi = U_\alpha \psi_{\varphi - \alpha}$. On the full $\dot{\mathcal{H}}$, however, we cannot in general define any \mathcal{U} .

This is the problem with gauge theory. When we quantize matter in the presence of background gauge fields, the resulting family of quantum theories in general realizes its classical gauge symmetry *via* a perfectly good ray representation. As far as the fermions are concerned there is *nothing wrong* with gauge symmetry. The phases in the ray representation are topologically unremovable; they prevent us from implementing the symmetry at all in the fully quantized theory, and in particular from imposing the constraint of gauge-invariance on the physical quantum states. Equivalently, in the temporal-gauge quantization of gauge theory [11, 12] we require that physical states obey

$$\left(\text{Tr} \left\{ T_\alpha \mathbf{D} \cdot \left(\frac{\delta}{\delta \mathbf{A}} \right) \right\} - i \psi^\dagger t_\alpha \psi \right) \Psi = 0, \quad (2)$$

which is the infinitesimal version of

$$\Psi[\mathbf{A}^g] = U_g \Psi[\mathbf{A}]. \quad (3)$$

But this just says that physical elements of $\dot{\mathcal{H}}$ must be *equivariant* sections of \mathcal{H} , or in other words that they must define sections of the reduced bundle $\bar{\mathcal{H}}$ over the true configuration space \mathcal{C} . If U_g is only projectively defined, then $\bar{\mathcal{H}}$ is not defined and this requirement makes no sense. If, moreover, the phases which spoil U_g have global topological content and so cannot be removed, then there is no cure for the problem. The theory is then anomalous.

A few remarks are in order before closing this section. We have established the existence of a nontrivial ray representation in a toy model by solving it exactly and noting the behavior of various eigenspaces of the energy globally over the parameter space. This brute-force approach will of course have to be replaced by something more powerful in field theory. Having established that at least one subbundle of \mathcal{H} twists on at least one orbit, we conclude that in the full theory the symmetry is “anomalous,” *i.e.*, it cannot be implemented as a true representation. Since the energy eigenspaces were all one-dimensional the only possible twist was the Möbius twist over a noncontractible circle in the symmetry group G . More generally we have to look for twists of higher-dimensional subbundles of \mathcal{H} , which will appear over higher-dimensional subspaces of G . In gauge theory, however, it will turn out to be enough to obtain a G -action on the vacuum subbundle, which is one-dimensional, and so there will be no anomalies due to obstructions beyond the first.

One might object that quantum mechanics involves not the real numbers but the complex, and that there are no interesting complex bundles over S^1 . We will answer this objection in two different ways in the sequel. For the $\pi_1(\mathcal{G}^3)$ anomaly, it is important for the G -action to preserve the real structure, while for the $\pi_2(\mathcal{G}^3)$ anomaly we indeed must consider two-spheres in \mathcal{G} (as the name implies). The former case resembles the obstruction to placing a spin structure on a space [13], since nontrivial $\pi_1(\mathcal{G}^3)$ implies nontrivial two-cells in \mathcal{C}^3 , while the latter resembles the obstruction to defining a spin^c structure, since it involves an *integer* (not Z_2) invariant and three-cells in \mathcal{C} .

3. Fermions

We begin for simplicity with the theory of [2], an $SU(2)$ gauge theory with a single isodoublet of Weyl fermions. This theory has a Euclidean Dirac operator which is strictly real [2]. Thus the energy eigenstates of the first-quantized theory can be chosen real, and the full second-quantized Hilbert bundle \mathcal{H} has a real structure.⁴ Furthermore, the representation matrix appearing in Gauss' law is real, and so the required \mathcal{G} -action must respect this real structure. As in our example, it will now suffice to show that the vacuum subbundle, say, is a Möbius bundle over any gauge orbit in order to establish the anomaly.

At each point of gauge configuration space we must now quantize fermions in the given background. This is not, of course, the usual procedure, in which one quantizes *free* fermions and treats gauge interactions perturbatively. Since the $SU(2)$ anomaly is nonperturbative, we must include the gauge fields from the start.

At the first-quantized level we encounter no difficulties. The Hilbert bundle is trivial, and the group action is $U_g v = v'$, where $v'(x) = g(x)v(x)$. Thus we expect any problems to come from second quantization, that is, from the definition of the Dirac sea. Accordingly let us focus our attention first on the vacuum subbundle \mathcal{H}_0 ; we will see that indeed once its \mathcal{G} -action has been defined there will be no further problems. Now the Dirac vacuum is defined as the state in Fock space in which all negative-energy states are filled. Since D^3 is gauge-covariant, all of its eigenvalues ϵ_i are gauge-invariant and \mathcal{H}_0 is mapped to itself by any gauge transformation (as indeed is any \mathcal{H}_ϵ filled to another Fermi level ϵ). Actually, though, \mathcal{H}_0 is unambiguously defined only on the subset \mathcal{A}' where none of the ϵ_i vanish. This turns out to be a small but crucial point, since unlike \mathcal{A} , which is contractible, \mathcal{A}' has nontrivial topology and so admits the possibility that the vacuum \mathcal{H}_0 can be twisted.

To establish the twist we combine the result of Berry [7], which relates twist to degeneracies, with the result of Witten [2], which establishes those degeneracies. Our argument is summarized in Fig. 1. Following Witten, we begin with the generator g^4 of $\pi_4(SU(2))$ and any point $A_{(0)}$ of \mathcal{A}^4 , the space of four-dimensional gauge potentials. Take $A_{(0)\mu}(x, t) \equiv 0$. Since \mathcal{A}^4 is connected, we can join $A_{(0)}$ to $[A_{(0)}]^{(g^4)}$ by a smooth path $A_{(\tau)}$, $\tau = 0$ to 1. For each τ we now transform $A_{(\tau)}$ by a time-dependent gauge transformation $g_{(\tau)}$ to put it into temporal gauge; call the result $A'_{(\tau)}$. In particular, $g_{(1)}$ is just $(g^4)^{-1}$, so instead of an open path of vector potentials each periodic in time we now have a closed loop of temporal-gauge histories, each of which ends at $A'_{(1)}(t = \infty) = [0]^{g_{(1)}(t = \infty)}$, a three-dimensional gauge transform of $A'_{(1)}(t = -\infty) \equiv 0$.

⁴ In particular, the vacuum subbundle \mathcal{H}_0 gets a real structure

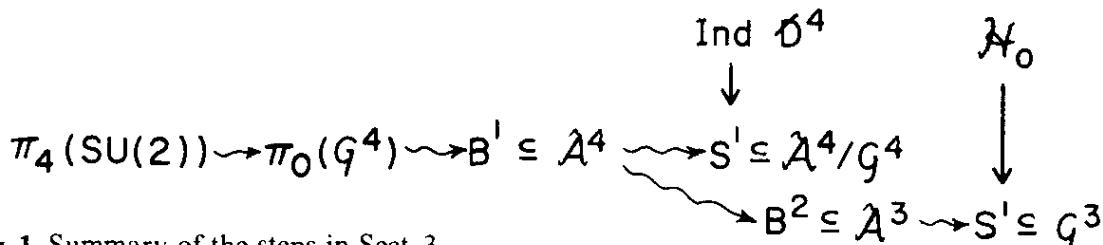


Fig. 1. Summary of the steps in Sect. 3

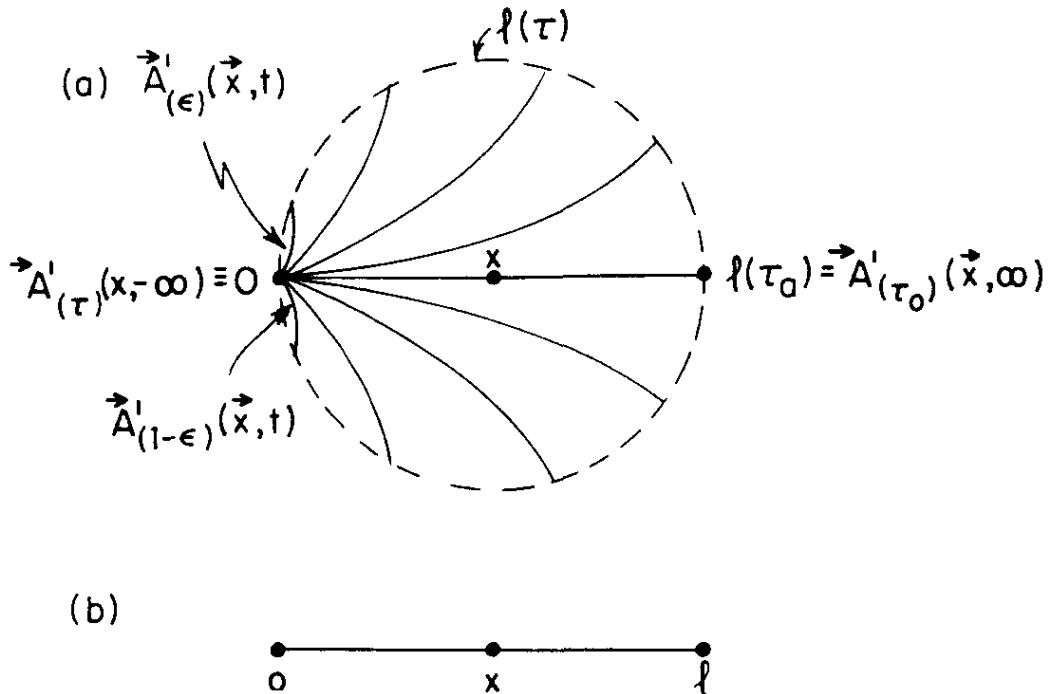


Fig. 2a and b. Disk in \mathcal{A}^3 associated to (a) SU(2) anomaly, (b) axial anomaly

The set of $A'_{(\tau)}(t)$, $-\infty < t < \infty$, $0 < \tau < 1$ thus forms a disk in \mathcal{A}^3 whose rim is a loop $\ell(\tau)$ of gauge transforms of zero (see Fig. 2). Each $\ell(\tau)$ is in \mathcal{A}' , since \mathcal{D}_0 has no zero modes on compactified space, and so we can restrict \mathcal{H}_0 to ℓ . We claim that $\mathcal{H}_0|_\ell$ is in fact twisted. For this to happen, there must be a point x on the disk excluded from \mathcal{A}' ; that is, there must be a degeneracy at x .

The presence of such a degeneracy follows at once from Witten's argument [2]. From the mod 2 index theorem, \mathcal{D}^4 must have a pair of zero modes at some $A_{(\tau_0)}$, and hence for the corresponding $A'_{(\tau_0)}$ as well. We can take these to be eigenstates ϕ_\pm of chirality. Taking Witten's argument one step further, if we choose each $A'_{(\tau)}$ to vary slowly in t then ϕ_\pm must be slowly-varying functions of time times eigenfunctions η_\pm^t of the Dirac hamiltonian $H_t \equiv \gamma_0 \mathcal{D}_{(\tau_0, t)}^3$. The energy eigenvalues must pass through zero⁵ at some t_0 , since ϕ_\pm are normalizable zero modes of the Euclidean \mathcal{D}^4 . Then $x = (\tau_0, t_0)$. Moreover, η_-^t has a CT-conjugated partner of opposite energy and chirality ζ_+^t , leading to the conical arrangement of left-handed energy eigenvalues shown in Fig. 3a. The number of these crossings will be equal, modulo two, to the number of Weyl isodoublets present.

When we second-quantize, the Fermi vacuum ray at each point $\ell(\tau)$ is the ray in Fock space with all negative-energy states filled. Choose a state $|0\rangle_0$ in this ray at 0. We can now attempt to adduce a nonvanishing section of $\mathcal{H}_0|_\ell$ by evolving $|0\rangle_0$ in the slowly-varying backgrounds $A'_{(\tau)}(t)$ for each τ . By the quantum adiabatic theorem [14], the final state will almost everywhere be almost pure vacuum, and we can project to \mathcal{H}_0 . This trick fails, however, at y . Here the adiabatic evolution passes through the vertex of the cone in Fig. 3a, producing the particle associated to η_+ and the antiparticle associated to ζ_+ . The resulting state has vanishing

⁵ Here is where the argument fails for the line bundle \mathcal{H}_ϵ filled up to some level other than zero, since the index theorem tells us nothing about ϵ -crossings

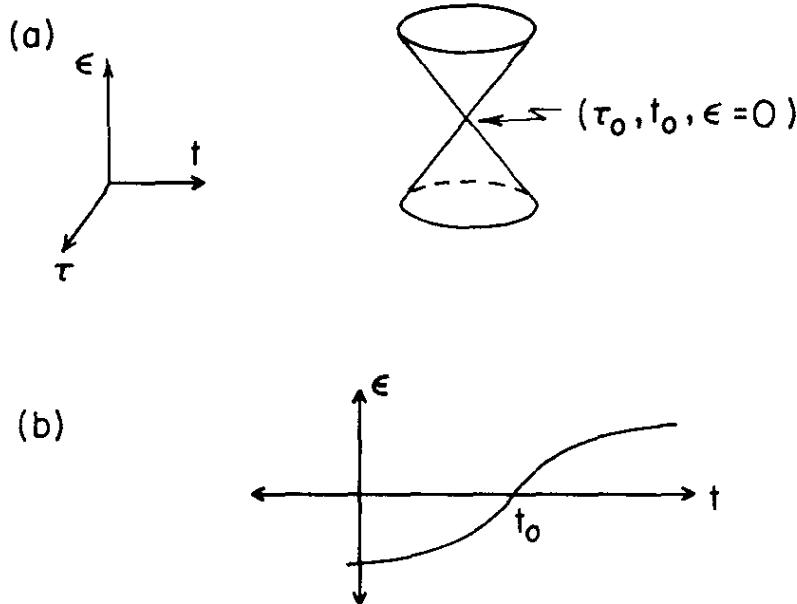


Fig. 3a and b. Eigenvalue behavior near $x = (r_0, t_0)$ for (a) SU(2) anomaly, (b) axial anomaly

projection to \mathcal{H}_0 . That is, the putative section “rolls over” near y , out of the plane of \mathcal{H}_0 and into an orthogonal direction provided by $\mathcal{H}_{\text{pair}}$. When projected to \mathcal{H}_0 , it vanishes at y . This reflects the twist of \mathcal{H}_0 .

This can all be made more precise using the result of [7]: For a loop of real hamiltonians, adiabatic transport around the loop returns to a state which is $(-1)^n$ times the original, if the loop encircles n simple degeneracies. (Note that the adiabatically-continued wavesections of this paragraph and one preceding were chosen only for convenience. Once we know that one section twists, we know they all do.)

While Berry’s result is elegant, we have given the pair-production picture as well in order to point up the physical similarities between the present case and the axial anomaly (see Figs. 2b and 3b). Usually the former is thought of in terms of phases, the latter in terms of particle production, but we can see that this really just a matter of emphasis. Particle production is crucial to *both*.⁶ In the case of the SU(2) anomaly, however, it occurs only for a special value τ_0 ; since it is not the generic behavior we do not find an important effect on the vacuum structure. Nevertheless, production is important, as it gives the sign twist which characterizes the anomaly. In the axial anomaly, on the other hand, it is production which is important in suppressing vacuum tunneling [11] while the phases do not matter. After all (Fig. 2b), in this case the rim of the disk is two *points* and so admits no twisted bundles.

Another important qualitative difference between the anomalies also comes from the codimension of x . In the SU(2) case, the level crossing had to be absent for points τ not exactly on τ_0 . For this to happen η^t_+ to have a partner ζ^t_+ , leading to zero net chirality production for the SU(2) anomaly. No such considerations apply

⁶ J. Goldstone has pointed out to us that our argument for the SU(2) anomaly is similar to one of his, summarized in [12], in which particle production also plays a key role

in the axial anomaly, and indeed (Fig. 3b) only η_+^i or ζ_+^i , not both, appears. Thus we get net production of chirality and a global symmetry is broken.

We can summarize the above discussion mathematically [15] by stating that the $\pi_0(\mathcal{G}^3)$ anomaly is given by the simplest invariant of the family index $\text{Ind } \mathbb{D}^4$ (a real virtual bundle over \mathcal{C}^4), namely its net dimension. Histories $A'(t)$ for which this is nonzero will have particle production with net change in chirality. The dimension is invariant to perturbations, so production is generic. The $\pi_1(\mathcal{G}^3)$ anomaly comes from the next invariant of $\text{Ind } \mathbb{D}^4$, its twist over circles in \mathcal{C}^4 . We have shown (Fig. 1) that this twist equals that of \mathcal{H}_0 over circles in \mathcal{A}^3 and so gives the obstruction to finding the \mathcal{G}^3 -action needed to quantize the theory correctly.⁷ For paths in \mathcal{A}^3 for which the lowest invariant vanishes, particle production gives no net chirality change and so comes from points where the null space of \mathbb{D}^4 jumps; i.e., production is not generic. No further invariants of \mathbb{D}^4 are relevant to \mathcal{H}_0 .

Unlike the $\pi_0(\mathcal{G}^3)$ anomaly, which is an integer, the $\pi_1(\mathcal{G}^3)$ anomaly can be cancelled by adding a second Weyl fermion of either chirality. Now a second pair state becomes degenerate with the vacuum and, by Berry's theorem, there is no sign change as we traverse ℓ .

We can also attempt to evade the anomaly by passing to the cover $\tilde{\mathcal{G}}^3$, as suggested in an earlier footnote. We now get a true $\tilde{\mathcal{G}}$ -action on \mathcal{H}_0 provided we map the nontrivial element \hat{g} covering the identity to the unitary operator -1 , and hence a $\tilde{\mathcal{G}}$ representation on \mathcal{H}_0 as in Sect. 2. If we take Gauss' law to mean that wave sections are equivariant under $\tilde{\mathcal{G}}$, however, we must in particular require that they be invariant under \hat{g} . Instead, all states have eigenvalue -1 under $\mathcal{U}_{\hat{g}}$! That is, we have succeeded only in defining on \mathcal{H} a ray realization of \mathcal{G} of the type studied in [6]. All states of \mathcal{H} are "fermionic." This construction recovers the formulation of the anomaly given in [2].

We have suggested that the anomaly is a second-quantization phenomenon, preventing us from finding an appropriate family of vacuum states. To go further, let us suppose that we have cancelled the obstruction and so have a well-defined \mathcal{G} -action on \mathcal{H}_0 . To get a \mathcal{G} -action on the rest of \mathcal{H} , we proceed as usual to define the Fock space creation operators a_A^{+i} on \mathcal{H}_A associated to the eigenfunctions η_A^i with energy $\varepsilon_A^i > 0$. (Similarly, b_A^i creates the mode η_A^i with energy $\varepsilon_A^i < 0$, and we reinterpret b_A^i as a destruction operator.) We can choose η_A^i smoothly in an open set V in \mathcal{A}^3 , and since there is an unambiguous \mathcal{G} -action on first-quantized states we can demand $\eta_{(AB)}^i(x) = g(x)\eta_A^i(x)$.

Now define

$$U_g(a_A^{+i_1} \dots a_A^{+i_k})|0\rangle_A \equiv a_{(AB)}^{+i_1} \dots a_{(AB)}^{+i_k} U_g|0\rangle_A, \quad (4)$$

for any vacuum state $|0\rangle_A$, and similarly with the b^+ . This definition is not arbitrary, but rather is dictated by the requirement that the quantum field built from a and b^+ have the same unambiguous transformation law as its first-quantized counterpart. Since there are no phase choices to make, there is no possibility of any obstruction to making them smoothly. Equation (4) defines a \mathcal{G} -action on a dense subspace of \mathcal{H}_A , $A \in V$. Furthermore, if on some other patch V_1 we choose a different

⁷ Since these twists are pure torsion, we cannot establish this fact by the use of real characteristic classes

orthonormal expansion η_{A1}^j (still equivariant under \mathcal{G}), we end up defining the same \mathcal{G} -action for \mathcal{H}_A , $A \in V \cap V_1$. Even as we approach a degenerate point, where \mathcal{H}_0 is not defined, we can extend this definition. Thus a \mathcal{G} -action on \mathcal{H}_0 extends without further difficulty to \mathcal{H} , and thence as we have seen to the full $\hat{\mathcal{H}}$.

We have therefore found that the higher invariants of $\text{Ind } \mathbb{D}^4$, like the lowest one, are irrelevant to implementing Gauss' law. All that matters is the twist of the index over circles. For the case to be discussed in Sect. 4, this agrees with the result of Atiyah and Singer [16], who use the path-integral formulation. It disagrees, however, with [15].

4. The Nonabelian Anomaly

The nonabelian anomaly presents almost no new features. An example of an affected theory is massless QCD with a triplet of left-handed Weyl "quarks." Since $\pi_5(\text{SU}(3)) = \pi_1(\mathcal{G}^4) = \mathbb{Z}$, we can consider the loop [17] in \mathcal{A}^4 given by transforming zero with each one of the noncontractible loop of $4d$ gauge transformations given by the generator g^5 . Again following the procedure outlined in Fig. 1, we then arrive at a three-ball in \mathcal{A}^3 whose boundary S^2 consists of three-dimensional gauge transformations of zero. Again by the family index theorem, D^4 generically has a pair of zero modes at one isolated value τ_0 , again leading to a conical vanishing of a pair of energy eigenvalues at some x in the interior of the ball. As we follow the trajectory given by τ_0 , we again find particle pair production obstructing the definition of a smooth nonvanishing vacuum section on the boundary of the ball. Berry's result for complex hamiltonians now says that indeed \mathcal{H}_0 is a twisted (monopole) line bundle over this S^2 ; its integer invariant is the nonabelian anomaly of the theory. Any action of \mathcal{G} now must have a string singularity somewhere, and so no acceptable version of Gauss' law exists.

Let us now attempt to pass to \mathcal{U}_g as before. Having established that \mathcal{H}_0 twists we can now forget about the interior of the orbit $\{\ell(\tau)\}$ and locally define our projective \mathcal{G} -action on $\mathcal{H}_0|_\ell$ as follows: Choose an S^2 metric on the orbit ℓ . If g is near the origin of \mathcal{G} and takes P to Q , $P, Q \in \ell$, consider the geodesic from P to Q as a slowly-varying history and evolve any vacuum state $|0\rangle_P$ in this background. Call the result $U_g|0\rangle_P \in \mathcal{H}_0|_Q$. Now suppose that h takes Q to R , also on the orbit; then $(hg)^{-1}$ takes R back to P . By a redefinition of phases we can now arrange for the adiabatic transport on the geodesic triangle so defined to return $|0\rangle_P$ multiplied by $e^{i\Omega/2}$, where Ω is the solid angle subtended by PQR [7]. The $\frac{1}{2}$ is fixed by the requirement that the phase factor be smoothly defined even for large g, h since then Ω is ambiguous by 4π ; this "Dirac quantization condition" on the normalization of the anomaly just reflects the fact that the anomaly is quantized due to its origin as a bundle twist.⁸

Thus $U_{hg}^{-1} U_h U_g |0\rangle_P = e^{i\Omega/2} |0\rangle_P$, and so $\mathcal{U}_{hg}^{-1} \mathcal{U}_h \mathcal{U}_g \Psi[P] = e^{i\Omega/2} \Psi[P]$. Choosing P to be any point where Ψ does not vanish we find once again that no state in $\hat{\mathcal{H}}$ is gauge-invariant.

Again we have seen that in the complex case the next-to-lowest invariant of the family index, in this case a two-form on \mathcal{C}^4 , is the only thing obstructing the

⁸ See also [17].

definition of a \mathcal{G}^3 -action on \mathcal{H} . Now, however, the obstruction is even more noticeable than in the previous case: since on a sphere we have nontrivial quantum holonomy even on infinitesimal loops,⁹ we expect that the $\pi_2(\mathcal{G}^3)$ anomaly should be visible even in perturbation theory. This is of course the case.

5. Remarks

There is an even more direct way to relate the lagrangian derivations of the anomaly to the hamiltonian picture. While it is less physical than the one given above, it does give the quickest way to find the *sign* of the integer invariant in the previous section, something we cannot do by examining the behavior of the energy eigenvalues alone. This sign was irrelevant in the Z_2 case; now we need it in order to recover the anomaly cancellation condition.

The lagrangian derivations show that the fermion partition function $e^{-\Gamma[A]}$ is actually a twisted section on \mathcal{C}^4 . In particular, it must vanish somewhere. But $e^{-\Gamma[A]}$ is just the vacuum expectation value of the time evolution operator $U(\infty, -\infty)$ in the presence of the time-dependent vector potential A_μ . Gauge transforming to temporal gauge as before, we get

$$\exp -\Gamma[A_{(\tau)}] =_{g_{(\tau)}} \langle 0 | U_{A'_{(\tau)}}(\infty, -\infty) | 0 \rangle_1. \quad (5)$$

Here $\{|0\rangle_g\}$ are a set of vacuum states on the various $\mathcal{H}_{(0)}$ s. Now as τ makes a complete circuit in the $\pi_1(\mathcal{G}^3)$ case, $A'_{(\tau)}$ returns to zero and so does its evolution operator. Since $e^{-\Gamma[A]}$ changes sign, it must be that the \mathcal{G} -action is twisted, as we found in Sect. 3. Furthermore, the single vanishing of $e^{-\Gamma[A]}$ which requires that it be twisted is just the signal of pair production again, since at τ_0 the evolved vacuum has no projection onto the transformed vacuum.

Repeating the argument in the case of the nonabelian $\pi_2(\mathcal{G}^3)$ anomaly, we find that not only must the \mathcal{G} -action be twisted, the twist in fact agrees in sign with that of the family index bundle. Hence the condition for the cancellation of the anomalous phases is that this bundle have no net twist, in agreement with [17]. In particular, ordinary QCD is safe.

From the hamiltonian point of view, the character of the gauge anomalies is determined by the structure of the possible real or complex line bundles over \mathcal{G}^3 . Loosely speaking, if over a gauge orbit \mathcal{H}_0 contains a unit of “flux” then it cannot be “squeezed” to zero, *i.e.*, the theory does not factor through to one properly defined on \mathcal{C}^3 . We have shown that the “flux” in a given theory’s configuration space can be computed solely in terms of the second invariant of its $\text{Ind } \mathcal{D}^4$. The fact that Witten’s anomaly appears only for symplectic groups like $SU(2)$, while the nonabelian anomaly appears for unitary groups like $SU(3)$ also comes naturally

⁹ This is codimension once again: on S^1 there are no interesting paths in a neighborhood of 0 which do not intersect 0. Note however that we do not claim to have obviated the perturbative analysis of gauge anomalies. As is well known, there are anomalies which the global analysis fails to uncover, either in the hamiltonian or lagrangian form. All we are saying is that when the global obstruction is present, it is clear why it makes its presence felt in perturbation theory.

from our construction, since in order to get interesting real (respectively complex) vacuum bundles over gauge orbits we needed nontrivial $\pi_1(\mathcal{G}^3)$ [respectively $\pi_2(\mathcal{G}^3)$]. This follows for the groups mentioned by the periodicity theorem.

While the higher invariants of the index are not related to gauge anomalies, they may still have interesting physical meaning, just as the lowest one does. The hamiltonian approach may yield further insight into this issue as well.

Note added. Some of the constructions in this paper have already been considered by I. M. Singer; see for example [19]. We thank the referee and Prof. Singer for bringing this work to our attention. After this paper was completed we also received the preprint by Faddeev [18], who discusses similar topics.

Acknowledgements. P. N. would like to thank O. Alvarez, S. Dellapietra, V. Dellapietra, J. Lott, N. Manton, and especially G. Moore for illuminating discussions, and the Institute for Theoretical Physics, Santa Barbara, for its hospitality while this work was being completed. We both thank R. Jackiw for useful discussions, and for bringing to our attention the preprint of Faddeev. This paper is based upon research supported in part by the National Science Foundation under Grant Nos. PHY77-27084 and PHY82-15249, supplemented by funds from the National Aeronautics and Space Administration.

References

1. Babelon, O., Viallet, C.: The Riemannian geometry of the configuration space of gauge theories. *Commun. Math. Phys.* **81**, 515 (1981)
2. Witten, E.: An SU(2) anomaly. *Phys. Lett.* **117B**, 324 (1982)
3. Bardeen, W.: Anomalous ward identities in spinor field theories. *Phys. Rev.* **184**, 1848 (1969)
Gross, D., Jackiw, R.: Effect of anomalies on quasi-renormalizable theories. *Phys. Rev. D* **6**, 477 (1972)
4. Coleman, S.: The uses of instantons. In: *The whys of subnuclear physics*. Zichichi, A. (ed.). New York: Plenum, 1979
Manton, N.: The Schwinger model and its axial anomaly. Santa Barbara preprint NSF-ITP-84-15 and references therein
5. Rajeev, S.: Fermions from bosons in $3+1d$ from anomalous commutators. *Phys. Rev. D* **29**, 2944 (1984)
6. Wigner, E.: Group theory New York: Academic, 1959; On unitary representations of the inhomogeneous Lorentz group. *Ann. Math.* **40**, 149 (1939)
7. Berry, M.: Quantal phase factors accompanying adiabatic changes. *Proc. R. Soc. Lond. A* **392**, 45 (1984)
8. Simon, B.: Holonomy, the quantum adiabatic theorem, and Berry's phase. *Phys. Rev. Lett.* **51**, 2167 (1983)
9. Wilczek, F., Zee, A.: Appearance of gauge structure in simple dynamical systems. *Phys. Rev. Lett.* **52**, 2111 (1984)
10. Atiyah, M.: *K-theory* New York: Benjamin 1967
11. Jackiw, R., Rebbi, C.: Vacuum periodicity in Yang-Mills theory. *Phys. Rev. Lett.* **37**, 172 (1977)
12. Jackiw, R.: Topological investigations of quantized gauge theories. Les Houches lectures, 1983 (MIT preprint CTP-1108)
13. Friedan, D., Windey, P.: Supersymmetric derivation of the Atiyah-Singer index and the chiral anomaly. *Nucl. Phys. B* **235** 395 (1984)
14. Messiah, A.: *Quantum mechanics*, Vol. 2. Amsterdam: North-Holland 1962
15. Sumitani, T.: Chiral anomalies and the generalized index theorem. Tokyo preprint UT-KOMABA84-7, 1984

16. Atiyah, M., Singer, I.: Dirac operators coupled to vector potentials. *Proc. Nat. Acad. Sci. USA* **81**, 2597 (1984)
17. Alvarez-Gaumé, L., Ginsparg, P.: The topological meaning of nonabelian anomalies. *Nucl. Phys.* **B243**, 449 (1984)
18. Faddeev, L.: Operator anomaly for gauss law. *Phys. Lett.* **145B**, 81 (1984). For this point of view on family ray representations, see also the recent preprints by R. Jackiw (MIT CTP 1209) and B. Zumino (Santa Barbara NSF-ITP-84-150)
19. Singer, I.: Families of dirac operators with applications to physics, M.I.T. Preprint, to appear in the proceedings of the Conference in Honor of E. Cartan, June 1984

Communicated by A. Jaffe

Received September 27, 1984; in revised form November 30, 1984

Chapter 8

CLASSICAL SYSTEMS

- [8.1] J. H. Hannay, "Angle Variable Holonomy in Adiabatic Excursion of an Integrable Hamiltonian," *J. Phys. A* **18** (1985) 221–230 426
- [8.2] M. V. Berry, "Classical Adiabatic Angles and Quantal Adiabatic Phase," *J. Phys. A* **18** (1985) 15–27 436
- [8.3] A. Shapere and F. Wilczek, "Gauge Kinematics of Deformable Bodies," to appear in *A. J. Phys.* 449
- [8.4] A. Shapere and F. Wilczek, "Geometry of Self-Propulsion at Low Reynolds Number," *J. Fluid Mech.* **198** (1989) 557–585 461

9

Asymptotics

It seems appropriate that we conclude with another beautiful contribution from Michael Berry. In this paper [9.1], he considers the construction of a systematic iterative approximation to the full phase, for which the geometric phase provides the leading term, in the original context of the spin 1/2 system. An asymptotic series is found, that has some quite remarkable and possibly universal properties. This work illustrates in a very concrete way how focusing on the concept of the geometric phase, has been a fruitful procedure.

One can imagine extending the analysis in several ways, particularly to treat degenerate levels. We believe, that many more attractive discoveries await determined explorers in these directions.

Quantum phase corrections from adiabatic iteration

By M. V. BERRY, F.R.S.

H. H. Wills Physics Laboratory, Tyndall Avenue, Bristol BS8 1TL, U.K.

(Received 22 April 1987)

The phase change γ acquired by a quantum state $|\psi(t)\rangle$ driven by a hamiltonian $H_0(t)$, which is taken slowly and smoothly round a cycle, is given by a sequence of approximants $\gamma^{(k)}$ obtained by a sequence of unitary transformations. The phase sequence is not a perturbation series in the adiabatic parameter ϵ because each $\gamma^{(k)}$ (except $\gamma^{(0)}$) contains ϵ to infinite order. For spin- $\frac{1}{2}$ systems the iteration can be described in terms of the geometry of parallel transport round loops C_k on the hamiltonian sphere. Non-adiabatic effects (transitions) must cause the sequence of $\gamma^{(k)}$ to diverge. For spin systems with analytic $H_0(t)$ this happens in a universal way: the loops C_k are sinusoidal spirals which shrink as ϵ^k until $k \sim \epsilon^{-1}$ and then grow as $k!$; the smallest loop has a size $\exp\{-1/\epsilon\}$, comparable with the non-adiabaticity.

1. INTRODUCTION

It is known (Berry 1984) that the phase of a quantum system whose hamiltonian is taken slowly round a cycle will acquire a geometric contribution, characteristic of the cycle, as well as the familiar dynamical one. The argument assumed that the instantaneous eigenstates were non-degenerate and that the adiabatic theorem could be applied. Two generalizations have removed these restrictions: Wilczek & Zee (1984) allow the instantaneous eigenstates to be degenerate (see also Segert 1987 and Mead 1987); and Aharonov & Anandan (1987) (see also Page 1987) allow the evolution to be non-adiabatic provided the system returns exactly to its initial state (apart from a phase, of course).

My purpose here is to develop a third generalization, going back to the original non-degenerate adiabatic scenario in which the system returns to its original state not exactly but in a close approximation, but now taking into account the finite rate at which the hamiltonian is changed. This leads to a technique for systematically obtaining corrections to the geometric phase. Garrison (1986) has made a start along these lines, by calculating the first phase correction in adiabatic perturbation theory. The method I shall use is not perturbative but iterative, and involves a sequence of unitary transformations chosen so as to make the hamiltonian cling ever more closely to the evolving state; it has the merit of being easy to visualize.

In §2 the phase is defined precisely and the iteration scheme described. It is shown how the phase can be interpreted as geometric or dynamical, depending on the choice of unitary transformation. The simplest non-trivial application (§3) is to a two-state (spin- $\frac{1}{2}$) system, for which the iteration can be formulated explicitly

in terms of geometry on the hamiltonian sphere. For finite slowness the evolution of the state will not be perfectly adiabatic, and as will be explained in §4 this implies the eventual divergence of the iteration scheme for the phase; for spin systems the divergence exhibits remarkable universality.

2. ITERATED ADIABATIC ANHOLONOMY

Let the state $|\psi_0(t)\rangle$ be driven by a hamiltonian $H_0(t)$ (the suffixes denote the zeroth state of the iteration scheme to be described below). In units with $\hbar = 1$, $|\psi_0\rangle$ satisfies

$$i|\dot{\psi}_0\rangle = H_0|\psi_0\rangle, \quad (1)$$

where here and hereafter time-dependences are understood where not written explicitly, and dots denote time derivatives. H_0 is taken smoothly round a cycle, i.e. $H_0(+\infty) = H_0(-\infty)$ with all derivatives vanishing as $|t| \rightarrow \infty$ (for this it suffices to take H_0 analytic in a strip including the real t axis). Let the (non-degenerate) instantaneous eigenstates of H_0 be $|n_0(t)\rangle$ with energies $E_0(n, t)$, i.e.

$$H_0|n_0\rangle = E_0(n)|n_0\rangle. \quad (2)$$

This defines the $|n_0\rangle$ up to a time-dependent phase which we make unique by demanding that

$$\langle n_0 | \dot{n}_0 \rangle = 0. \quad (3)$$

With this choice, the eigenstates are parallel-transported, as explained by Simon (1983). Let the system start in the n th eigenstate, i.e.

$$|\psi_0(-\infty)\rangle = |n_0(-\infty)\rangle \equiv |N\rangle. \quad (4)$$

The phase which is the object of study is now defined as

$$\gamma(n) \equiv \text{Im} \ln \langle N | \psi_0(+\infty) \rangle + \int_{-\infty}^{\infty} dt E_0(n, t). \quad (5)$$

This form of writing assumes that the integral over E_0 (which is minus the dynamical phase) converges, or can be made to converge by shifting the energy origin; if not, γ can be defined by a suitable limiting procedure. As defined by (5) the phase is more general than that which arises in the cyclic evolutions of Aharonov & Anandan (1987), because transitions may (and usually do) make $|\langle N | \psi_0(+\infty) \rangle| < 1$, i.e. the final state may be a superposition including states other than the original. Here, however, the emphasis is on cases where such non-adiabatic effects are small. We shall introduce a slowness parameter ϵ , entering H_0 in the combination ϵt , and regard ϵ as small. In the limit $\epsilon = 0$, γ becomes the geometric phase studied previously (Berry 1984).

Now let us follow several other authors (e.g. Avron *et al.* 1987; Anandan & Stodolsky 1987; Mead 1987) and define $U_0(t)$ as the unitary operator generating the eigenstates $|n_0(t)\rangle$ by acting on the original eigenstates $|N\rangle$, i.e.

$$|n_0(t)\rangle = U_0(t)|N\rangle. \quad (6)$$

Quantum phase corrections from adiabatic iteration

33

Because of the parallel-transport law (3), $|n_0(+\infty)\rangle$ differs from $|N\rangle$ by a phase which is precisely the original geometric phase $\gamma_0(n)$ (anholonomy of H_0), so that

$$U_0(+\infty)|N\rangle = \exp\{i\gamma_0(n)\}|N\rangle. \quad (7)$$

The operator U_0 naturally leads to a new representation of the evolving state $|\psi_0\rangle$ as that state $|\psi_1\rangle$ on which it must act to produce $|\psi_0\rangle$, i.e.

$$|\psi_1\rangle = U_0^\dagger|\psi_0\rangle. \quad (8)$$

Thus

$$\langle N|\psi_0(+\infty)\rangle = \langle N|U_0(+\infty)|\psi_1(+\infty)\rangle = \exp\{i\gamma_0\}\langle N|\psi_1(+\infty)\rangle \quad (9)$$

so that the phase (5) now becomes

$$\gamma(n) = \gamma_0(n) + \text{Im} \ln \langle N|\psi_1(+\infty)\rangle + \int_{-\infty}^{\infty} dt E_0(n, t). \quad (10)$$

To proceed further, we need the Schrödinger equation satisfied by $|\psi_1\rangle$. This involves a hamiltonian $H_1(t)$ which differs from H_0 , because the transformation U_0 is time-dependent. Thus

$$i|\dot{\psi}_1\rangle = H_1|\psi_1\rangle, \quad (11)$$

where

$$H_1 = U_0^\dagger H U_0 - i U_0^\dagger \dot{U}_0. \quad (12)$$

It is not difficult to show from (2), (3) and (6) that the matrix elements of H_1 in the $|N\rangle$ representation are

$$\langle M|H_1|N\rangle = E_0(n)\delta_{MN} - \frac{i\langle m_0|\dot{H}_0|n_0\rangle}{E_0(n)-E_0(m)}(1-\delta_{MN}). \quad (13)$$

The simplest adiabatic approximation is to neglect the off-diagonal elements on the grounds that \dot{H}_0 is of order ϵ . Then (11) gives

$$|\psi_1(t)\rangle \approx \exp\left\{-i\int_{-\infty}^t dt' E_0(n, t')\right\}|N\rangle \quad (14)$$

so that (10) gives

$$\gamma(n) \approx \gamma_0(n). \quad (15)$$

Systematic improvements can, however, be achieved by not neglecting the off-diagonal terms. Instead, $H_1(t)$ is regarded as a new hamiltonian with new eigenstates $|n_1(t)\rangle$ and new eigenvalues $E_1(n, t)$, and the transformation repeated, leading to a new representation $|\psi_2(t)\rangle$ and a further hamiltonian $H_2(t)$. Obviously the procedure can be iterated according to the scheme

$$\begin{aligned} H_{k+1} &= U_k^\dagger H_k U_k - i U_k^\dagger \dot{U}_k, \\ \text{i.e. } \langle M|H_{k+1}|N\rangle &= E_k(n)\delta_{MN} - \frac{i\langle M_k|\dot{H}_k|n_k\rangle}{E_k(n)-E_k(m)}(1-\delta_{MN}), \end{aligned} \quad (16)$$

where $U_k|N\rangle = |n_k\rangle$, $H_k|n_k\rangle = E_k(n)|n_k\rangle$, $\langle n_k|n_k\rangle = 0$.

(Iteration does not change the initial states $|N\rangle$, because of the assumed smoothness of H_0 .)

At the k th iteration step,

$$\begin{aligned}\langle N | \psi_0(+\infty) \rangle &= \langle N | U_0(+\infty) U_1(+\infty) \dots U_k(+\infty) | \psi_{k+1}(+\infty) \rangle \\ &= \exp \left\{ i \sum_{j=0}^k \gamma_j(n) \right\} \langle N | \psi_{k+1}(+\infty) \rangle,\end{aligned}\quad (18)$$

where $\gamma_j(n)$ are the anholonomies of the cycled hamiltonians $H_j(t)$. We can stop at this iteration by neglecting the off-diagonal elements in H_{k+1} . This gives the k th phase approximant

$$\gamma(n) \approx \gamma^{(k)}(n) \equiv \sum_{j=0}^k \gamma_j(n) + \int_{-\infty}^{\infty} dt [E_0(n, t) - E_k(n, t)]. \quad (19)$$

The sequence of approximants is not a perturbation series in the adiabatic parameter ϵ , because even $\gamma^{(1)}$ involves ϵ to infinitely high order (of course $\gamma^{(0)} = \gamma_0$ is independent of ϵ). Rather, the iterations can be regarded as successive superadiabatic transformations to moving frames (in Hilbert space) attempting to cling ever more closely to the evolving state $|\psi_0\rangle$ (we will see in §4 that the attempts ultimately fail).

It is instructive to digress and consider iteration schemes that are not based on the parallel-transport law (3). An obvious class of alternatives (infinitely many) is to require the $|n_k(t)\rangle$ to return exactly to $|N\rangle$ as $t \rightarrow +\infty$. This would eliminate the anholonomies of the U_k , but would change the diagonal elements of the iterated hamiltonians to

$$\langle N | H_{k+1}(t) | N \rangle = E_k(n) - i \langle n_k | \dot{n}_k \rangle, \quad (20)$$

where E_k and $|n_k\rangle$ are of course different from those in the previous scheme. Instead of (19) we would have

$$\gamma(n) \approx \int_{-\infty}^{\infty} dt [E_0(n, t) - E_k(n, t) - \text{Im} \langle n_k | \dot{n}_k \rangle], \quad (21)$$

so that in this form of iteration the phase arises from the approximate eigenvalue of H_{k+1} , and so its derivation appears entirely ‘dynamical’, even for $k = 0$! (γ_0 itself is geometric, regardless of how it is derived, because it does not depend on ϵ). Obviously it is possible to construct intermediate iteration schemes, in which the derivation appears as partly anholonomic and partly dynamical. The precise classical analogue of this interpretational ambiguity can be seen in the contrasting treatments of adiabatic angles by Hannay (1985) (anholonomic) and Berry (1985) (dynamical). The reason for choosing the iteration scheme based on (3) is firstly that it is unique and secondly that the corrections to the hamiltonian at each stage (16) are entirely in the off-diagonal terms and thus higher order in ϵ .

3. SPIN- $\frac{1}{2}$ SYSTEMS

We take

$$H_0(t) = \mathbf{R}_0(t) \cdot \boldsymbol{\sigma}, \quad (22)$$

where $\mathbf{R}_0 \equiv (X_0, Y_0, Z_0)$ and $\boldsymbol{\sigma}$ is the vector spin- $\frac{1}{2}$ operator. Thus

$$H_0 = \frac{1}{2} \begin{pmatrix} Z_0 & X_0 - i Y_0 \\ X_0 + i Y_0 & -Z_0 \end{pmatrix}. \quad (23)$$

Quantum phase corrections from adiabatic iteration

35

The transformation (12) will generate a new hamiltonian $H_1(t) = \mathbf{R}_1(t) \cdot \boldsymbol{\sigma}$, and the aim of this section is to find the explicit form of this operator.

For simplicity of writing we temporarily omit the suffixes zero. Define unit vectors $\mathbf{r}(t)$, $\mathbf{v}(t)$, $\mathbf{w}(t)$, and positive scalars $R(t)$, $V(t)$, by

$$\mathbf{R} \equiv R\mathbf{r}, \quad \dot{\mathbf{r}} \equiv V\mathbf{v}, \quad \mathbf{w} \equiv \mathbf{r} \times \mathbf{v}, \quad (24)$$

and think of \mathbf{r} as a radius vector of the unit sphere. Then the cycle of $H_0(t)$ is represented by transport of the triad $\mathbf{r}, \mathbf{v}, \mathbf{w}$ round a circuit C_0 on the sphere (figure 1). The eigenvalues of H are $\pm \frac{1}{2}\mathbf{R}(t)$, and the corresponding eigenstates $|n(t)\rangle$ will be denoted by $|\pm(t)\rangle$, or often simply by $|\pm\rangle$; these states depend on $\mathbf{r}(t)$ but not $R(t)$.

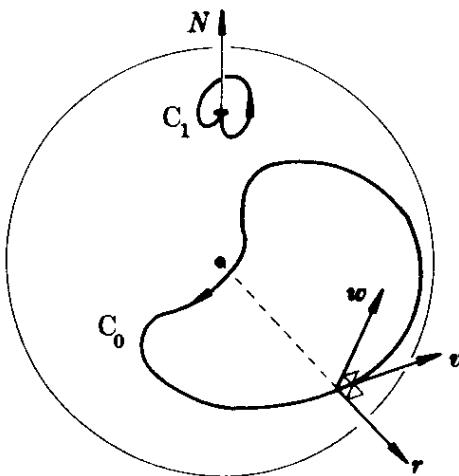


FIGURE 1. Transport of triad $\mathbf{r}, \mathbf{v}, \mathbf{w}$ round initial circuit C_0 on the hamiltonian sphere, and iterated circuit C_1 .

The operator $U(t)$ (6) turns the eigenstates $|\pm(-\infty)\rangle$ into the eigenstates $|\pm(t)\rangle$ by a sequence of infinitesimal rotations. It is shown in appendix A that the instantaneous angular velocity of the rotation is uniquely determined by (3), that is $\langle \pm | \dot{\pm} \rangle = 0$, to be the angular velocity $\Omega_{||}(t)$ of a frame *parallel-transported* with \mathbf{r} over the spheres. This angular velocity is

$$\Omega_{||} = \mathbf{r} \times \dot{\mathbf{r}} = V\mathbf{w} \quad (25)$$

and must be distinguished from the different angular velocity Ω_t of the triad $\mathbf{r}, \mathbf{v}, \mathbf{w}$, which is

$$\Omega_t = \mathbf{v} \times \dot{\mathbf{v}}. \quad (26)$$

Thus $U(t)$ is the time-ordered product

$$U(t) = T \exp \left\{ -i \int_{-\infty}^t dt' \Omega_{||}(t') \cdot \boldsymbol{\sigma} \right\}. \quad (27)$$

To find the new operator H_1 , we now use (12), which gives

$$H_1 = U^\dagger (\mathbf{R} - \Omega_{||}) \cdot \boldsymbol{\sigma} U. \quad (28)$$

This result has a purely classical origin (valid for any spin) in the transformation, to a non-inertial frame moving with \mathbf{R} and rotating with angular velocity $\Omega_{||}$, of

the equation of motion for the expectation value $\langle \sigma \rangle$. (See, for example, Cina (1986), Suter *et al.* (1987), Anandan & Stodolsky (1987), and, for explicit calculations in terms of Hannay's angle, Berry (1986) and Gozzi & Thacker (1987)). From (13) and (28), the matrix representation in terms of the initial states is

$$H_1 = \begin{pmatrix} \frac{1}{2}R & iV\langle +|\sigma \cdot \mathbf{v}|-\rangle \\ -iV\langle -|\sigma \cdot \mathbf{v}|+\rangle & -\frac{1}{2}R \end{pmatrix}, \quad (29)$$

where use has been made of the fact that the off-diagonal elements of $\sigma \cdot \mathbf{r}$ vanish.

To make this explicit we must evaluate $\langle +|\sigma \cdot \mathbf{v}|-\rangle$. In Appendix A this is shown to be

$$i\langle +|\sigma \cdot \mathbf{v}|-\rangle = \frac{1}{2} \exp \left\{ -i \int_{-\infty}^t dt' \Omega(t') \right\}, \quad (30)$$

where

$$\Omega \equiv (\Omega_t - \Omega_{||}) \cdot \mathbf{r} = \mathbf{v} \times \dot{\mathbf{v}} \cdot \mathbf{r} = \mathbf{w} \cdot \dot{\mathbf{v}}. \quad (31)$$

Thus Ω measures the rate at which the \mathbf{rvw} triad twists about \mathbf{r} , relative to the parallel-transported frame. The new hamiltonian (29) now becomes $H_1(t) = \mathbf{R}_1(t) \cdot \sigma$ where (reinstating the suffixes)

$$Z_1 = R_0, \quad X_1 + iY_1 = V_0 \exp \left\{ i \int_{-\infty}^t dt' \Omega_0(t') \right\}, \quad (32)$$

As $\mathbf{r}_0(t)$ executes the loop C_0 representing H_0 , the new unit vector executes a loop C_1 representing H_1 (figure 1). Because V_0 is adiabatically small (it is the speed at which \mathbf{r}_0 moves on the unit sphere), this new loop is very close to the north pole of the sphere. In fact C_1 resembles a cardioid with a cusp or corner at the pole, because $V_0 \rightarrow 0$ as $|t| \rightarrow \infty$. Iteration of the map from \mathbf{R}_0 to \mathbf{R}_1 gives $\mathbf{R}_2(t)$, $\mathbf{R}_3(t)$... and hence further loops C_2 , C_3 , ... For small ϵ these loops rapidly diminish in size at first, and we might hope that they will continue to do so (especially because $R_{k+1} = (R_k^2 + V_k^2)^{\frac{1}{2}} > R_k$ so iteration takes H further from the degeneracy at $R = 0$). But this hope cannot be realized, as will be explained in §4.

To obtain the phase approximations $\gamma^{(k)}(\pm)$ from (19) we must determine the anholonomies $\gamma_j(\pm)$ from the unitary operators (27). These can be obtained by noting that the effect of all the infinitesimal rotations in the product (27) from $t = -\infty$ to $t = +\infty$ is a spinor rotation about the initial direction $\mathbf{r}(-\infty)$ (which of course is the same as the final direction) by the parallel transport angles A_j associated with the loops C_j . These are simply the *solid angles* subtended by the loops at the origin of the sphere, obtained from (31) as the total twist of the parallel frame about the \mathbf{rvw} triad, plus the 2π rotation of that triad, namely

$$\begin{aligned} A_j &= 2\pi - \int_{-\infty}^{\infty} dt \Omega_j(t) = 2\pi - \int_{-\infty}^{\infty} dt \mathbf{v} \times \dot{\mathbf{v}} \cdot \mathbf{r} \\ &= 2\pi - \oint_{C_j} d\mathbf{r} \cdot \mathbf{r}' \times \mathbf{r}, \end{aligned} \quad (33)$$

where primes denote differentiation with respect to arc length on the sphere.

Thus with the z axis temporarily along $\mathbf{r}(\pm\infty)$ we have

$$U_j(+\infty) = \exp\left\{-\frac{1}{2}iA_j\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}\right\} = 1 \cos(\frac{1}{2}A_j) - i\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \sin(\frac{1}{2}A_j) \quad (34)$$

so that

$$U_j(+\infty) \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \exp\{-\frac{1}{2}iA_j\} a \\ \exp\{+\frac{1}{2}iA_j\} b \end{pmatrix} \quad (35)$$

giving

$$\gamma_j(\pm) = \mp\frac{1}{2}A_j, \quad (36)$$

as found by Berry (1984). The energies in (19) are simply $E_k(\pm) = \pm\frac{1}{2}R_k$, so the phase approximants are

$$\gamma^{(k)}(\pm) = \mp\frac{1}{2} \left\{ \sum_{j=0}^k A_j + \int_{-\infty}^{\infty} dt [R_k(t) - R_0(t)] \right\}. \quad (37)$$

As illustration, consider the lowest-order approximations to the phase when C_0 is a circle of latitude on the unit sphere, with polar angle θ and azimuth $\phi(t)$ satisfying $\phi(+\infty) - \phi(-\infty) = 2\pi$. Thus

$$\mathbf{R}_0 = \mathbf{r}_0 = (\sin\theta \cos\phi(t), \sin\theta \sin\phi(t), \cos\theta). \quad (38)$$

The speed (24) on the sphere is

$$V_0 = \sin\theta\dot{\phi} \quad (39)$$

and the twist (31) is

$$\Omega_0 = \cos\theta\dot{\phi} \quad (40)$$

so (32) gives the new loop C_1 as $\mathbf{r}_1(t) = \mathbf{R}_1(t)/R_1(t)$, where

$$\mathbf{R}_1(t) = (\sin\theta\dot{\phi} \cos\{\phi(t)\cos\theta\}, \sin\theta\dot{\phi} \sin\{\phi(t)\cos\theta\}, 1). \quad (41)$$

Now we introduce an adiabatic ϵ by

$$\phi(t) \equiv \Phi(\tau), \quad \tau \equiv \epsilon t. \quad (42)$$

Then it is not difficult to show that the anholonomies A_j in (37) are even in ϵ while the ‘dynamical’ integrals over R_k are odd. Apart from

$$A_0 = 2\pi(1 - \cos\theta), \quad (43)$$

which is of course independent of ϵ , all the other terms in (37) are of infinite order in ϵ . Thus

$$\int_{-\infty}^{\infty} dt [R_1(t) - R_0(t)] = \frac{1}{\epsilon} \int_{-\infty}^{\infty} d\tau [(1 + \epsilon^2 \sin^2\theta \Phi'^2(\tau))^{\frac{1}{2}} - 1], \quad (44)$$

where the prime denotes $d/d\tau$. The first three terms of an expansion are contained in the first iteration $\gamma^{(1)}(\pm)$, and are

$$\gamma_{\pm}^{(1)}(\pm) \approx \mp \left\{ \pi(1 - \cos\theta) + \frac{1}{4}\epsilon \sin^2\theta \int_{-\infty}^{\infty} d\tau \Phi'^2(\tau) + \frac{1}{4}\epsilon^2 \sin^2\theta \cos\theta \int_{-\infty}^{\infty} d\tau \Phi'^3(\tau) \right\}. \quad (45)$$

The terms originate in A_0 , R_1 and A_1 , respectively. Garrison (1986) has obtained the first two terms of (45) by adiabatic perturbation (ϵ -expansion).

4. INEVITABILITY OF DIVERGENCE

To get the k th phase approximant (19) we neglected the off-diagonal terms in H_{k+1} and so approximated $\langle N | \psi_{k+1}(+\infty) \rangle$ in (18) by a pure phase factor. This ignores transitions to other states, which will cause the survival probability $|\langle N | \psi_{k+1}(+\infty) \rangle|^2$ to deviate from unity. The transition probability is typically exponentially small (Hwang & Pechukas 1977), i.e.

$$\Delta(\epsilon) \equiv 1 - |\langle N | \psi_k(+\infty) \rangle|^2 \sim \exp\{-1/\epsilon\}; \quad (46)$$

$\Delta(\epsilon)$ is independent of the order of iteration k (cf. (18)). The sequence $\gamma^{(k)}(n)$ cannot converge for finite ϵ because this would imply $\Delta = 0$. Therefore the sequence must diverge: the true phase must reflect the non-analyticity of the survival amplitude, unlike the terms in (19), which although of infinite order in ϵ are nevertheless analytic at $\epsilon = 0$.

We expect the terms $\gamma_j(n)$ to get smaller at first and then increase. This is the typical behaviour of an asymptotic expansion (Dingle 1973), and it is reasonable to hope that the best approximant is the one for which $|\gamma^{(k+1)}(n) - \gamma^{(k)}(n)|$ is smallest (excluding perversities such as oscillations in this quantity) and that this value is of the same order as the non-adiabaticity (see Balian *et al.* (1978) for a numerical demonstration of a related phenomenon).

We can illustrate the inevitable divergence with the loops C_j on the unit sphere (figure 1) which represents the successive hamiltonians $H_j(t)$ for a spin- $\frac{1}{2}$ particle. For small ϵ the first iterated loop C_1 , and many subsequent ones, will be small and close to the north pole. If we write the radius vector as $r = (x, y, z)$ then $z \approx 1$ and we can approximate the loops as lying in the tangent plane at the pole. Parallel transport is now ordinary euclidean parallel translation, and in the iteration (32) $\int_{-\infty}^t dt' \Omega(t')$ is the direction of the tangent to the loop at t , relative to that at $t = -\infty$. Denoting $x + iy$ by ζ we find that (32) reduces to the simple iteration

$$\zeta_{j+1}(t) = \dot{\zeta}_j(t)/R_j(t). \quad (47)$$

When investigating the behaviour of the loops thus generated we can set $R_j(t) = 1$, because the successive radii differ little from $R_0(t)$ (cf. 31 which shows that $R_1^2 = R_0^2 + V_0^2$) and can be reduced to unity by the time rescaling $dt \rightarrow R_0(t) dt$. Thus (47) becomes the iterated hodograph transformation of mechanics, namely

$$\zeta_k(t) = d^k \zeta_0(t)/dt^{k+1}. \quad (48)$$

The loop C_k is generated in the xy plane by letting t run from $-\infty$ to $+\infty$. To describe adiabatic circuits we consider t to appear in the combination et . Then ζ_k contains a factor e^k and the C_k initially decrease in size. It would be reasonable to expect this decrease to continue, but it does not. The surprising fact is that, for almost all initial loops $\zeta_0(t)$, the C_k ultimately get bigger, and moreover for small ϵ the nature of the increase is *universal*. Furthermore, the winding number of C_k , defined as the number of rotations of the tangent as the loop is traversed, exceeds that of C_{k-1} by $\frac{1}{2}$ (alternate loops have cusps at $\zeta = 0$).

To justify these assertions we begin by recalling the assumption that $H_0(t)$, and

hence $\zeta_0(t)$ is smooth, and interpret this as analyticity in a strip about the t axis. Thus the Fourier transform $\bar{\zeta}_0(\omega)$, defined by

$$\zeta_0(t) = \int_{-\infty}^{\infty} d\omega \bar{\zeta}_0(\omega) \exp\{-i\omega et\} \quad (49)$$

decays exponentially as $|\omega| \rightarrow \infty$, the exponents being the imaginary parts of the singularities of $\zeta_0(t)$ nearest to the real axis in the upper and lower halves of the et plane. If these singularities are

$$\tau_+ = \tau_{1+} + i\tau_{2+}, \quad \tau_- = \tau_{1-} - i\tau_{2-} \quad (\tau_{2+}, \tau_{2-} > 0) \quad (50)$$

$$\text{then } \bar{\zeta}_0(\omega) \rightarrow A_{\pm} \exp\{i\omega\tau_{1\pm}\} \exp\{-|\omega|\tau_{2\pm}\} \quad \text{as } \omega \rightarrow \pm\infty. \quad (51)$$

The iterated loops (49) are given by (48) as

$$\zeta_k(t) = (-i)^k \int_{-\infty}^{\infty} d\omega (\epsilon\omega)^k \bar{\zeta}_0(\omega) \exp\{-i\omega et\}. \quad (52)$$

For large k only the asymptotic form (51) of $\bar{\zeta}_0(\omega)$ contributes (because of the $(\epsilon\omega)^k$ factor) so that

$$\zeta_k(t) \rightarrow (-i\epsilon)^k k! \{A_+/[\tau_{2+} - i(\tau_{1+} - \epsilon t)]^{k+1} + (-1)^k A_-/[\tau_{2-} + i(\tau_{1-} - \epsilon t)]^{k+1}\}. \quad (53)$$

The term with the smaller of $\tau_{2\pm}$ dominates exponentially, so that after a trivial shift of time origin and redefinition of ϵ as $\epsilon/\min\{\tau_{2\pm}\}$ the k th loop takes the universal form

$$\begin{aligned} \zeta_k(t) \rightarrow & \frac{A(i\epsilon)^k k!}{(1-i\epsilon t)^{k+1}} = A[\epsilon^k k!/(1+\epsilon^2 t^2)^{\frac{1}{2}(k+1)}] \\ & \times \exp\{i[\frac{1}{2}k\pi + (k+1)\arctan\epsilon t]\} \quad \text{as } k \rightarrow \infty. \end{aligned} \quad (54)$$

These universal loops C_k are Maclaurin's sinusoidal spirals (Lawrence 1972), some of which are shown in figure 2. In polar coordinates defined by $r = i^k \exp\{i\phi\}$ their equation is

$$r_k(\phi) = A\epsilon^k k! \cos^{k+1}\{\phi/(k+1)\} \quad (55)$$

(C_0 is a circle, C_1 a cardioid and C_2 Cayley's sextic). The maximum radius (at $\phi = 0$) is $r_k = A\epsilon^k k!$. This decreases at first (because of ϵ^k) but ultimately increases (because of $k!$). The smallest maximum radius occurs when $k \approx \epsilon^{-1}$ and is

$$r_{1/\epsilon}(0) \approx A(2\pi/\epsilon)^{\frac{1}{2}} \exp\{-1/\epsilon\}. \quad (56)$$

The k th loop has winding number $\frac{1}{2}k+1$; the initial and final windings ($t \approx \pm(k+1)/\pi\epsilon$) have radii smaller than the largest radius ($t = 0$) by a factor $[\pi/(k+1)]^{k+1}$.

The universality of the sinusoidal spirals can be described alternatively by saying that these curves are the attractors of the hodograph map in the space of loops. In view of the well-known instability of differentiation, the existence of attractors is remarkable, especially when considered backwards: the almost spirals C_k (k large), when iterated under the inverse map which is an integration and therefore supposedly stabilizing, must diversify into the infinite variety of

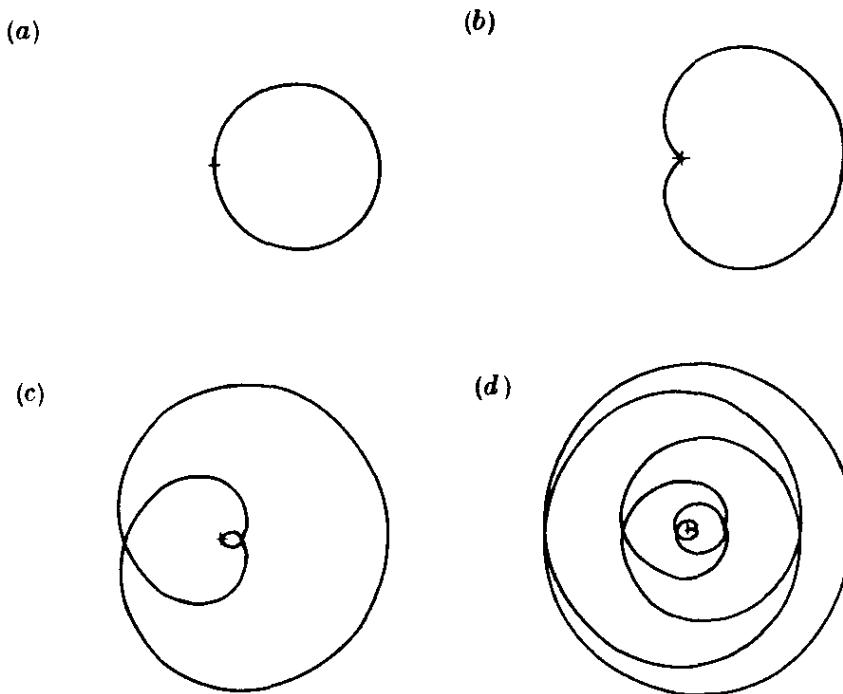


FIGURE 2. Universal loops (55) (sinusoidal spirals) for (a) $k = 0$; (b) $k = 1$; (c) $k = 9$; (d) $k = 50$. The loops are normalized to have the same maximum distance from the origin (which is reached at $t = \pm\infty$). The k th loop has $\frac{1}{2}k + 1$ windings (not all visible in (c) and (d) because they are so small).

possible C_0 . The resolution of the apparent paradox must lie in the assumed analyticity of C_0 .

The asymptotic behaviour of the loops is reflected in the phase approximants $\gamma^{(k)}$. In the north-pole plane approximation (48) the solid angle A_k is just the euclidean area of C_k , namely

$$A_k = \frac{1}{2} \int_{-\infty}^{\infty} dt \operatorname{Im} \zeta_k^* \dot{\zeta}_k \quad (57)$$

and the difference between successive radii is (cf. 32)

$$R_{k+1} - R_k = (R_k^2 + V_k^2)^{\frac{1}{2}} - R_k \approx \frac{1}{2} |\dot{\zeta}_k|^2 \quad (58)$$

because we have normalized R_k approximately to unity. Thus, roughly, the difference between successive approximants (37) is, with (48),

$$|\gamma^{(k+1)} - \gamma^{(k)}| \approx \frac{1}{4} \int_{-\infty}^{\infty} dt \{ \operatorname{Im} (\dot{\zeta}_k^* \dot{\zeta}_k) + |\dot{\zeta}_k|^2 \} \quad (59)$$

With the universal (54) this becomes

$$|\gamma^{(k+1)} - \gamma^{(k)}| \approx \frac{1}{4} A^2 [(k+1)!]^2 \epsilon^{2k+1} (\epsilon(k+\frac{3}{2}) + 1) I_{k+2}, \quad (60)$$

where $I_k \equiv \int_{-\infty}^{\infty} d\tau / (1 + \tau^2)^k = 2 \int_0^{\frac{1}{2}\pi} d\theta \cos^{2(k-1)} \theta \rightarrow \left(\frac{\pi}{k}\right)^{\frac{1}{2}}$ as $k \rightarrow \infty$. (61)

For small ϵ , (60) falls to a minimum when $k \approx \epsilon^{-1}$, whose value is of order $\exp\{-1/\epsilon\}$. When units are reinstated, the critical k acquires the following physical interpretation: ϵ^{-1} is the ratio of the average transition frequency between the two instantaneous eigenstates and the average rate of adiabatic change of

$\ln H_0(t)$. With this interpretation, $\exp\{-1/\epsilon\}$ is indeed the size of the non-adiabaticity (Hwang & Pechukas 1977).

It is possible to construct loops C_0 whose iterates do not fall into the universality class just discussed. One way is to make $\zeta_0(t)$ have its ‘nearest singularity’ actually *on* the real t axis whilst possessing all derivatives. For example, if $\zeta_0 \sim 1 - \exp\{-1/\epsilon|t|\}$ the smallest loop occurs for iteration $k \sim \epsilon^{-\frac{1}{2}}$ and has radius of order $\exp\{\epsilon^{-\frac{1}{2}}\}$. The opposite situation is for $\zeta_0(t)$ to be so smooth that its ‘nearest singularity’ lies infinitely far from the axis. For example, if $\zeta_0 \sim \exp\{-(\epsilon t)^2\}$ the smallest loop occurs for iteration $k \sim \epsilon^{-2}$ and has radius of order $\exp\{\epsilon^{-2}\}$. Another possibility is for the imaginary parts $\tau_{2\pm}$ in (50) to be equal, so that neither of the two contrary loops in (53) dominates; one example is when $\zeta_0(t)$ is real, that is C_0 is simply a back-and-forth swinging enclosing no area (then there is no anholonomy and γ is purely dynamical). It does seem, however, that whatever the form of $\zeta_0(t)$ (provided all derivatives exist, and vanish as $|t| \rightarrow \infty$) the loops do eventually grow, but I have not been able to construct a general proof.

The nonadiabaticity, as expressed by the transition probability $\Delta(\epsilon)$ defined by (46), can be estimated by perturbation theory, using $|N\rangle$ as the unperturbed state. For spin- $\frac{1}{2}$ systems it is tempting to employ perturbation theory in the north pole plane approximation, applying it to the hamiltonian for the iteration for which C_k is smallest. However, this gives only a crude approximation, for reasons worth exploring because they illuminate a curious general feature of the adiabatic approximation.

Application of standard time-dependent perturbation theory to the iterated hamiltonian (13) gives the transition probability $\Delta(\epsilon)$ in (46) as the sum of the probabilities of transitions to states $M \neq N$:

$$\Delta(\epsilon) \approx \sum_{M \neq N} \left| \int_{-\infty}^{\infty} dt \frac{\langle m_0 | \dot{H}_0 | n_0 \rangle}{E_0(n) - E_0(m)} \exp \left\{ i \int_{-\infty}^t dt' [E_0(m) - E_0(n)] \right\} \right|^2, \quad (62)$$

For a spin- $\frac{1}{2}$ system, starting (say) in $|+\rangle$, this gives, for the probability that at $t = +\infty$ there has been a transition to $|-\rangle$,

$$\Delta(\epsilon) \approx \frac{1}{4} \left| \int_{-\infty}^{\infty} \frac{V_0(t)}{R_0(t)} \exp \left\{ i \int_{-\infty}^t dt' [\Omega_0(t') - R_0(t')] \right\} \right|^2, \quad (63)$$

where use has been made of (29) and (30). Both the preceding formulae are independent of the order of iteration.

In the north-pole plane approximation with $R_0 = 1$, (63) becomes

$$\Delta(\epsilon) \approx \frac{1}{4} \left| \int_{-\infty}^{\infty} dt \zeta_0(t) \exp\{-it\} \right|^2. \quad (64)$$

This too is independent of the order of iteration, as can be seen from (48) and integration by parts. Substituting the universal loop (54) gives

$$\begin{aligned} \Delta(\epsilon) &\approx \frac{A^2}{4} \left| \int_{-\infty}^{\infty} dt (1 - i\epsilon t)^{-1} \exp\{-it\} \right|^2 \\ &= \frac{A^2 \pi^2}{\epsilon^2} \exp\{-2/\epsilon\} \quad \text{if } \epsilon > 0 \\ &= 0 \quad \text{if } \epsilon < 0 \end{aligned} \quad (65)$$

The two cases $\epsilon > 0$ and $\epsilon < 0$ correspond to opposite senses for the traversals of the loops and hence to an original hamiltonian $H_0(t)$ and its time reverse $H_0(-t)$. It is not surprising that time reversal can lead to very different transition probabilities because almost all the hamiltonians which give rise to anholonomy lack time-reversal symmetry. What is surprising – and this is the curious feature mentioned earlier – is that the phase is insensitive to this qualitative distinction to all orders of adiabatic iteration: the approximants (19) merely change sign under time reversal.

However, the extreme difference between $\epsilon > 0$ and $\epsilon < 0$ in (65) is an artefact of the north-pole plane approximation, arising from the fact that for $\epsilon < 0$ the perturbation $(1 - i\epsilon t)^{-1}$ has no negative-frequency components to stimulate the transition from $|+\rangle$ to $|-\rangle$. Without this approximation, but still taking C_0 as a circle, not necessarily small, as given by (38), the perturbation formula (63) gives, on making use of (39) and (40),

$$\Delta(\epsilon) \approx \frac{\sin^2 \theta}{4} \left| \int_{-\infty}^{\infty} dt \dot{\phi}(t) \exp \{i[\phi(t) \cos \theta - t]\} \right|^2. \quad (66)$$

Now take

$$\phi(t) = 2 \arctan \epsilon t, \quad (67)$$

which gives the same form of cycling as the universal loop (54). Then integration by parts gives

$$\Delta(\epsilon) \approx \frac{\tan^2 \theta}{4\epsilon^2} \left| \int_{-\infty}^{\infty} d\tau \left(\frac{1+i\tau}{1-i\tau} \right)^{\cos \theta} \exp \{-i\tau/\epsilon\} \right|^2. \quad (68)$$

Changing the sign of ϵ has the same effect as changing θ to $\pi - \theta$, as it should because the effects on $|+\rangle$ of time reversal and latitude reversal of C_0 are the same.

In appendix B it is shown that the asymptotic form of (68) for small positive ϵ is

$$\Delta(\epsilon) \approx [\sin \pi \cos \theta) \tan \theta \Gamma(1 - \cos \theta) 2^{\cos \theta}]^2 \exp \{-2/\epsilon\} / \epsilon^{2 \cos \theta}. \quad (69)$$

Near the poles and the equator this has the limiting forms

$$\begin{aligned} \Delta(\epsilon) &\rightarrow 4\pi^2 \theta^2 \exp \{-2/\epsilon\} / \epsilon^2 && \text{as } \theta \rightarrow 0, \\ &\rightarrow \pi^2 \exp \{-2/\epsilon\} && \text{as } \theta \rightarrow \frac{1}{2}\pi, \\ &\rightarrow \frac{1}{16}\pi^2 (\pi - \theta)^6 \epsilon^2 \exp \{-2/\epsilon\} && \text{as } \theta \rightarrow \pi. \end{aligned} \quad (70)$$

The sense in which (65) is a crude approximation is now evident; instead of a discontinuity between $\epsilon > 0$ and $\epsilon < 0$ there is a smooth transition involving θ , which shows that time reversal does not make Δ vanish but reduces its value by ϵ^4 .

5. CONCLUDING REMARKS

The main result of this work is the formula (19) giving the phase approximants obtained from a (unique) succession of unitary transformations clinging ever closer to the evolving state. Successive approximants $\gamma^{(k)}(n)$ are correct to higher orders in the adiabatic parameter ϵ , but (19) is not a power series because each

approximant (except the lowest) contains ϵ to infinite order. Nevertheless, at all orders of iteration the scheme neglects non-adiabatic transitions, and because these are of order $\exp\{-1/\epsilon\}$ this quantity, rather than ϵ itself, can be regarded as the adiabatic parameter and then the entire sequence of approximants can be considered to be contained within the lowest-order adiabatic approximation (higher approximations to the exact phase (5) would involve powers of $\exp\{-1/\epsilon\}$).

For spin systems the unitary sequence can be interpreted geometrically as (initially) shrinking loops on the hamiltonian sphere. The loops could be observed, for example, by exploiting the fact that the expectation $\langle \sigma \rangle$ evolves classically according to $\langle \sigma \rangle = \mathbf{R}_0 \times \langle \sigma \rangle$. Then the unitary sequence corresponds to a sequence of transformations to rotating frames. Stopping at the k th such transformation, making the adiabatic approximation and transforming back to the original reference frame, we find that $\langle \sigma \rangle$ follows not $\mathbf{R}_0(t)$ but $\mathbf{R}^{(k)}(t)$ which (cf. 28) includes corrections from the angular velocities $\boldsymbol{\Omega}_j(t)$ of the successive frames, i.e.

$$\mathbf{R}^{(k)}(t) = \mathbf{R}_0(t) - \sum_{j=0}^{k-1} \boldsymbol{\Omega}_{jj}(t). \quad (71)$$

Now, successive $\boldsymbol{\Omega}_j$ are (initially) smaller (by ϵ), and generate loops C_k (§4) which in their own frames have increasing winding numbers. Therefore the motion of $\langle \sigma \rangle$ is a sequence of ever-finer nutations forming a hierarchy reminiscent of Ptolemy's epicycles. Successive orders of iteration correspond to observing the motion with increasing resolution.

Several questions are raised by the divergence of the sequence of phase approximants. One concerns the universality of the shrinking-and-growing of the hamiltonian loops C_k for spin systems, when k is large, ϵ small and $H_0(t)$ analytic. This was derived within the north-pole plane approximation, and slight doubt lingers as to whether the result would survive inclusion of the effects of the curvature of the sphere. Assuming it does, another question is whether the divergence for non-spin systems has a similar adiabatic universality. (Of course for spins the universality we are discussing can last only for iteration numbers not much greater than $k \approx \epsilon^{-1}$; subsequent iterations will take the expanded loops C_k away from the north pole plane. What happens then? Is there an infinite sequence – possibly irregular – of further shrinking and growing when $k \gg \epsilon^{-1}$? Or do the windings continue to increase, leading to ultimate loops C_∞ covering the sphere densely?)

Finally, one wonders whether there are any systems for which the iteration would converge (or stop at some order) because the adiabatic approximation would be exact. This question is prompted by a spatial analogy: the existence of one-dimensional potentials $V(x)$ for which the semiclassical ('adiabatic') approximation is exact for some energies E , so that there is no reflection (i.e. no 'transition'). For example, with arbitrary real 'quantum momentum' $k_{qu}(x)$ the local plane wave

$$\psi(x) = \exp \left\{ i \int_0^x dx' k_{qu}(x') \right\} / [k_{qu}(x)]^{\frac{1}{2}} \quad (72)$$

44

M. V. Berry

is an exact solution of

$$\psi''(x) + k_{\text{cl}}^2(x) \psi(x) = 0 \quad (73)$$

provided it is related to the classical momentum $k_{\text{cl}} = [(E - V)]^{\frac{1}{2}}$ by

$$k_{\text{qu}} = [k_{\text{cl}}^2 + k_{\text{qu}}^{\frac{1}{2}}(k_{\text{qu}}^{-\frac{1}{2}})'']^{\frac{1}{2}}. \quad (74)$$

Then there is no coupling to the reflected wave which is the complex conjugate of (72). We can think of k_{qu} as the outcome of infinite order semiclassical iteration of (74), the lowest (W.K.B.) approximation being $k_{\text{qu}} \approx k_{\text{cl}}$. The analogue of γ is the phase of the transmitted wave (in this case entirely dynamical) referred to the W.K.B. phase, i.e.

$$\begin{aligned} \gamma &= \int_{-\infty}^{\infty} dx (k_{\text{qu}} - k_{\text{cl}}) \\ &= \int_{-\infty}^{\infty} dx k_{\text{qu}}^{\frac{1}{2}} (k_{\text{qu}}^{-\frac{1}{2}})'' / \{k_{\text{qu}} + [k_{\text{qu}}^2 - k_{\text{qu}}^{\frac{1}{2}} (k_{\text{qu}}^{-\frac{1}{2}})'']^{\frac{1}{2}}\}, \end{aligned} \quad (75)$$

which from (74) is of infinite order in spatial slowness. The analogous adiabatic problem, of constructing a state $|\psi(t)\rangle$ evolving exactly as an eigenstate of some changing $H(t)$, seems much more difficult. Successive iterated hamiltonians would have to commute with each other. This is impossible for non-trivial spin systems, but it is conceivable that in other cases there would occur a conspiracy of the off-diagonal elements $\langle m_k | \dot{H}_k | n_k \rangle$ to allow at least one of the states to evolve adiabatically. (Of course I am not here denying the existence of the cyclic evolutions considered by Aharonov & Anandan 1987, because these involve states that return to themselves without having at every instant to be eigenstates of the driving hamiltonian; the spatial analogy is the Ramsauer-Townsend effect, in which perfect transmission is achieved without requiring $\psi(x)$ to be everywhere locally plane as in (72).)

APPENDIX A. SPIN- $\frac{1}{2}$ CALCULATIONS JUSTIFYING (27) AND (30)

The defining equation for $U(t)$ is (6), which can be written

$$U(t)|\pm(-\infty)\rangle = |\pm(t)\rangle. \quad (A 1)$$

Differentiating leads to

$$\dot{U}U^\dagger|\pm\rangle = |\dot{\pm}\rangle, \quad (A 2)$$

which on making use of $|+\rangle\langle+| + |-\rangle\langle-| = 1$ becomes

$$\dot{U} = (|\dot{+}\rangle\langle+| + |\dot{-}\rangle\langle-|)U \equiv BU. \quad (A 3)$$

Then (27) follows if

$$B = -iV\mathbf{w} \cdot \boldsymbol{\sigma}. \quad (A 4)$$

Quantum phase corrections from adiabatic iteration

45

The diagonal elements $\langle \pm |B| \pm \rangle$ vanish because of orthogonality and $\langle \pm | \dot{\pm} \rangle = 0$. The diagonal elements $\langle \pm |w \cdot \sigma| \pm \rangle$ also vanish because $|\pm\rangle$ are eigenstates of $r \cdot \sigma$, and w is perpendicular to r . The off-diagonal elements are

$$\langle \pm |B| \mp \rangle = \langle \pm | \dot{\mp} \rangle. \quad (\text{A } 5)$$

Differentiating the eigenequations

$$r \cdot \sigma | \pm \rangle = \pm \frac{1}{2} | \pm \rangle \quad (\text{A } 6)$$

gives

$$\langle \pm | \dot{\mp} \rangle = \mp V \langle \pm | w \cdot \sigma | \mp \rangle. \quad (\text{A } 7)$$

Now choose local axes in which r, v, w lie along z, x, y respectively. Then

$$\langle + | \dot{-} \rangle = -\frac{V}{2} (1 \ 0) \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = -\frac{V}{2} \quad (\text{A } 8)$$

and $\langle + |B| \rangle = -iV \langle + | w \cdot \sigma | - \rangle = -i \frac{V}{2} (1 \ 0) \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \frac{-V}{2}. \quad (\text{A } 9)$

Arguing similarly for $\langle - | \dot{+} \rangle$ we see that the operator equation (A 4) holds for all matrix elements and hence is true, thus implying (27).

To prove (30) we first show that $|\langle + | \sigma \cdot v | - \rangle| = \frac{1}{2}$:

$$\begin{aligned} \langle + | \sigma \cdot v | - \rangle \langle - | \sigma \cdot v | + \rangle &= \langle + | \sigma \cdot v | (-) \langle - | + | + \rangle \langle + | \sigma \cdot v | + \rangle \\ &= \langle + | (\sigma \cdot v)^2 | + \rangle = \frac{1}{4}. \end{aligned} \quad (\text{A } 10)$$

(The first equality is valid because $\langle + | \sigma \cdot v | + \rangle = 0$.) Thus we can write

$$i \langle + | \sigma \cdot v | - \rangle = \frac{1}{2} \exp \{i\mu(t)\} \quad (\text{A } 11)$$

To find an equation for μ , differentiate:

$$\frac{1}{2} \dot{\mu} \exp \{i\mu\} = \langle + | \dot{\sigma} \cdot v | - \rangle + \langle + | \sigma \cdot \dot{v} | - \rangle + \langle + | \sigma \cdot v | \dot{-} \rangle. \quad (\text{A } 12)$$

The first and last terms vanish because

$$\langle \dot{\pm} | \sigma \cdot v | \mp \rangle = \langle \dot{\pm} | \pm \rangle \langle \pm | \sigma \cdot v | \mp \rangle + \langle \dot{\pm} | \mp \rangle \langle \mp | \sigma \cdot v | \mp \rangle = 0. \quad (\text{A } 13)$$

In the middle term, v can be replaced by its component along w , which from (31) gives, with (A 11)

$$\dot{\mu} = -i\Omega \frac{\langle + | \sigma \cdot w | - \rangle}{\langle + | \sigma \cdot v | - \rangle}. \quad (\text{A } 14)$$

The matrix elements are just those previously evaluated in (A 7–A 9), so that

$$\dot{\mu} = -\Omega, \quad (\text{A } 15)$$

which with (A 10) and (A 11) gives (30).

APPENDIX B. ASYMPTOTIC EVALUATION OF THE INTEGRAL (68)

Because of the exponent in (68), the integration contour can be deformed (for positive ϵ) into the negative half-plane to surround a cut extending from the branch point at $\tau = -i$ to $\tau = -i\infty$. On the right side we can take $\tau =$

$-i + r \exp\{-\frac{1}{2}i\pi\}$ (r from 0 to ∞) and on the left we can take $\tau = -i + r \exp\{\frac{3}{2}i\pi\}$ (r from ∞ to 0). Thus the integral becomes

$$\begin{aligned} & -i \exp\{-1/\epsilon\} \int_0^\infty dr \frac{\exp\{-r/\epsilon\} (2+r)^{\cos\theta}}{r^{\cos\theta}} [\exp(i\pi \cos\theta) - \exp(-i\pi \cos\theta)] \\ &= \frac{-i \exp\{-1/\epsilon\}}{\epsilon^{(\cos\theta-1)}} \int_0^\infty dx \frac{\exp\{-x\} (2+\epsilon x)^{\cos\theta}}{x^{\cos\theta}} 2i \sin(\pi \cos\theta). \end{aligned}$$

For small ϵ the term ϵx can be neglected. The resulting integral is a Γ -function and squaring gives (69).

REFERENCES

- Aharonov, Y. & Anandan, J. 1987 *Phys. Rev. Lett.* **58**, 1593–1596.
 Anandan, J. & Stodolsky, L. 1987 *Phys. Rev. D* **35**, 2597–2600.
 Avron, J. E., Seiler, R. & Yaffe, L. G. 1987 *Communs math. Phys.* **110**, 33–49.
 Balian, R., Parisi, G. & Voros, A. 1978 *Phys. Rev. Lett.* **41**, 1141–1144.
 Berry, M. V. 1984 *Proc. R. Soc. Lond. A* **392**, 45–57.
 Berry, M. V. 1985 *J. Phys. A* **18**, 15–27.
 Berry, M. V. 1986 *Adiabatic phase shifts for neutrons and photons*. In *Fundamental aspects of quantum theory* (ed. V. Gorini & A. Frigerio). NATO ASI series vol. 144, pp. 267–278. New York: Plenum.
 Cina, J. 1986 *Chem. Phys. Lett.* **132**, 393–395.
 Dingle, R. B. 1973 *Asymptotic expansions: their derivation and interpretation*. New York and London: Academic Press.
 Garrison, J. C. 1986 Preprint UCRL 94267, Lawrence Livermore Laboratory.
 Gozzi, L. E. & Thacker, W. D. 1987 *Phys. Rev. D* **35**, 2388–2398.
 Hannay, J. H. 1985 *J. Phys. A* **18**, 221–230.
 Hwang, J.-T. & Pechukas, P. 1977 *J. Chem. Phys.* **67**, 4640–4653.
 Lawrence, J. D. 1972 *A catalog of special plane curves*. Dover Publications.
 Mead, C. A. 1987 *Phys. Rev. Lett.* **59**, 161–164.
 Page, D. H. 1987 *Phys. Rev. Lett.* (Submitted.)
 Segert, J. 1987 *J. math. Phys.* (In the press.)
 Simon, B. 1983 *Phys. Rev. Lett.* **51**, 2167–2170.
 Suter, D., Chingas, G., Harris, R. A. & Pines, A. 1987 *Molec. Phys.* (In the press.)
 Wilczek, F. & Zee, A. 1984 *Phys. Rev. Lett.* **52**, 2111–2114.