



OXFORD JOURNALS
OXFORD UNIVERSITY PRESS

The British Society for the Philosophy of Science

Understanding Electromagnetism

Author(s): Gordon Belot

Reviewed work(s):

Source: *The British Journal for the Philosophy of Science*, Vol. 49, No. 4 (Dec., 1998), pp. 531-555

Published by: [Oxford University Press](#) on behalf of [The British Society for the Philosophy of Science](#)

Stable URL: <http://www.jstor.org/stable/688130>

Accessed: 09/07/2012 06:48

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <http://www.jstor.org/page/info/about/policies/terms.jsp>

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



Oxford University Press and *The British Society for the Philosophy of Science* are collaborating with JSTOR to digitize, preserve and extend access to *The British Journal for the Philosophy of Science*.

<http://www.jstor.org>

Understanding Electromagnetism

Gordon Belot

ABSTRACT

It is often said that the Aharonov–Bohm effect shows that the vector potential enjoys more ontological significance than we previously realized. But how can a quantum-mechanical effect teach us something about the interpretation of Maxwell’s theory—let alone about the ontological structure of the world—when both theories are false? I present a rational reconstruction of the interpretative repercussions of the Aharonov–Bohm effect, and suggest some morals for our conception of the interpretative enterprise.

1 Introduction

2 Gauge theories and their interpretation

3 Interpreting electromagnetism

3.1 Interpretations

4 Quantization and the Aharonov–Bohm effect

4.1 Quantizations and interpretation

4.1 Quantizing the charged particle in a magnetic field

4.3 The Aharonov–Bohm effect

5 Conclusion

1 Introduction

When one first learns classical electromagnetism, one is taught to think of Maxwell’s equations as governing the evolution in time of the electric and magnetic fields (or, more subtly, of the electromagnetic field). Under this interpretation the theory is both deterministic and local: *deterministic* in the sense that specifying the present state of the fields suffices to fix their past and future; *local* in the sense that if we want to know what will happen next *here*, the theory tells us that we need only look at the field values *hereabouts* right now—we do not need to know what is happening arbitrarily far away. Thus construed, electromagnetism is the paradigm of all that a classical (i.e. non-quantum) theory should be.

Although this way of thinking about electromagnetism remains the pedagogical standard, it has been known for some time to be untenable. In 1959, Aharonov and Bohm argued that a charged *quantum* particle moving in the region external to a solenoid would be sensitive to whether or not current were running through the device, *despite the fact that the field values in the regions*

of space occupied by the particle would be unaffected by the operation of the solenoid.

It is widely agreed that the subsequent experimental detection of the Aharonov–Bohm effect discredited the familiar way of understanding electromagnetism. One can maintain the traditional interpretation of the theory only by maintaining that fields act where they are not. But this flies in the face of the well-entrenched principle that classical fields act by contact rather than at a distance. It would seem, then, that the electric and magnetic fields cannot constitute the ontology of electromagnetism. It is now standard to maintain that the Aharonov–Bohm effect shows that the vector potential, formerly viewed as a mere mathematical convenience, must in fact be physically real. The advantage being that in the region exterior to the solenoid the vector potential—unlike the magnetic field—depends on the state of the device, so that one can explain the behaviour of the particle in terms of the values of vector potential in the region it actually occupies.

Thus, it is often said that the Aharonov–Bohm effect shows that the traditional interpretation of electromagnetism must be replaced. I subscribe to this conclusion. But I would put it somewhat differently: *until the discovery of the Aharonov–Bohm effect, we misunderstood what electromagnetism was telling us about our world.* This formulation captures what I take to be the kernel of the common wisdom. But it is intentionally provocative: it brings to the fore the epistemological and metaphysical puzzles inherent in episodes like the post-Aharonov–Bohm reinterpretation of electromagnetism.

After all, by the time the Aharonov–Bohm effect was discovered, it had long since been accepted that electromagnetism does not accurately describe our world. In an influential paper of 1933, Bohr and Rosenfeld argued that there can be no consistent theory of the interaction between charged quantum particles and a classical electromagnetic field.¹ Thus our world could not possibly contain the sort of field described by Maxwell's equations: electromagnetism is a false theory. Now there is a very straightforward sense in which a false—but eminently useful—theory like electromagnetism can tell us about our world: it makes empirical predictions which are very accurate within a certain circumscribed domain of applicability. But it seems strange to say that the *interpretation* of such a theory tells us about our world. To interpret a theory is to describe the possible worlds about which it is the literal truth. Thus, an interpretation of electromagnetism seems to tell us about a possible world which we *know* to be distinct from our own. On the other hand, whatever world

1 Bohr and Rosenfeld's paper is reprinted as Ch. 1 of Cohen and Stachel [1979]. The models which are used to predict the Aharonov–Bohm effect are idealizations in which the classical field acts on the particles, but the particles are not sources of the field. Thus, although they may be useful for describing certain phenomena, they cannot be taken to be accurate representations of our world any more than can the models of electromagnetism.

electromagnetism is true of, it is not one which contains quantum electrons. So it is difficult to see how a quantum-mechanical effect can teach us anything about the interpretation of electromagnetism. Of course, quantum mechanics itself is false (being nonrelativistic). So *our* world is one about which *neither* electromagnetism *nor* quantum mechanics is true. None the less, I maintain, we learn something about our own world when we study the interpretative interaction between these two false theories.

I will present a logical reconstruction of the Aharonov–Bohm effect which suggests a resolution of the tension between falsehood and interpretative interest. My first task is to describe and situate the formalism and interpretative problems of electromagnetism. I begin in Section 2 by sketching the formalism and interpretative problems of gauge theories in general. This allows me to present electromagnetism as a gauge theory in Section 3, and to show how its interpretative problems arise out of its gauge freedom. The presentation of these two sections is quite abstract, but I believe that this approach gives valuable insight into the structure of the classical theory. In Section 4, I show how attention to quantum mechanics can shift the balance of power among competing interpretations of electromagnetism: there is an ambiguity inherent in the construction of a quantum model of a charged particle moving in an electromagnetic field; distinct interpretations of electromagnetism suggest different ways of resolving this ambiguity; the empirical success of one or another quantum treatment can then have repercussions for our attempts to interpret electromagnetism. In particular, we will see that in the aftermath of the Aharonov–Bohm effect, we are forced to accept that electromagnetism is either indeterministic or nonlocal. Thus we find that the requirement that our false theories mesh in an appropriate way—ontologically as well as empirically—places strong constraints upon our interpretative practice. In the final section of the paper, I attempt to explicate a sense in which false theories tell us about our world, and to show how this fact has important consequences for our understanding of the structure and content of our physical knowledge.

Before beginning my main task, I would like to say a few words about a distinction which will play a fundamental role in what follows. I distinguish three components of a physical theory: the formalism, the interpretation, and the application. The formalism is some (more or less rigorous) mathematics. This might be of interest to a mathematician with no interest whatsoever in physics. The application is a set of practices which allow one to derive and to test the empirical consequences of the theory. The interpretation consists of a set of stipulations which pick out a putative ontology for the possible worlds correctly described by the theory. Schematically, we can imagine the physical theory being taught in a course for undergraduates: the formalism is developed on the blackboard during lectures; the application is worked out in problem sets and in the lab; the interpretation is fixed via verbal asides which give the

students a heuristic grasp of the content of the theory. A command of all three components will be essential for any student who aspires to full understanding of the theory.

Now, of course, in saying that we misunderstood electromagnetism prior to the discovery of the Aharonov–Bohm effect, I do not mean to suggest that Einstein misunderstood the formalism of the theory or that Hertz misunderstood its application. These remain fixed as our interpretation changes (or, more properly, they evolve via their own dynamics which need not be directly correlated with interpretative developments). But I *do* insist that a change in interpretation constitutes a change in understanding.

2 Gauge theories and their interpretation

I am going to sketch a couple of frameworks for doing classical mechanics, and discuss their respective interpretative problems.² One is the familiar Hamiltonian formalism. The other is a generalization of the Hamiltonian framework: the language of gauge theories. They share a great deal of their conceptual apparatus. Both Hamiltonian systems and gauge systems consist of triples of mathematical objects: a space, a tensor which gives this space some geometric structure, and a real-valued function on the space, called the *Hamiltonian*. When equipped with its geometric structure, the space is called the *phase space*, and its points are thought of as representing the dynamically possible states of some classical physical system (typically a set of particles or fields). The Hamiltonian then determines a class of curves in phase space. These are thought of as representing the dynamically possible histories of the system—if we know which point represents the present state of the system, then a curve through this point passes through points representing the dynamically possible future and past states of the system.

This much, Hamiltonian systems and gauge systems have in common. The difference between them lies in the nature of the geometric structure of phase space. As we will see, the weaker geometric structure of gauge systems brings with it a thorny interpretative problem.

(i) Hamiltonian systems

The geometric structure of the phase space of a Hamiltonian system is called a *symplectic form*.³ Its chief virtue is the following: specifying a real-valued function on the phase space (the Hamiltonian, H) suffices to determine a unique *dynamical trajectory* through each point of phase space.⁴ Figure 1 is

² The following presentation is meant to be accessible to a general reader. Details are given in the footnotes.

³ The phase space consists of a manifold, M , equipped with a closed, nondegenerate, two-form, ω .

⁴ These curves are the integral curves of the vector field X_H which solves $X_H \lrcorner \omega = dH$ (the left-hand side of this equation is the contraction of X_H with ω).

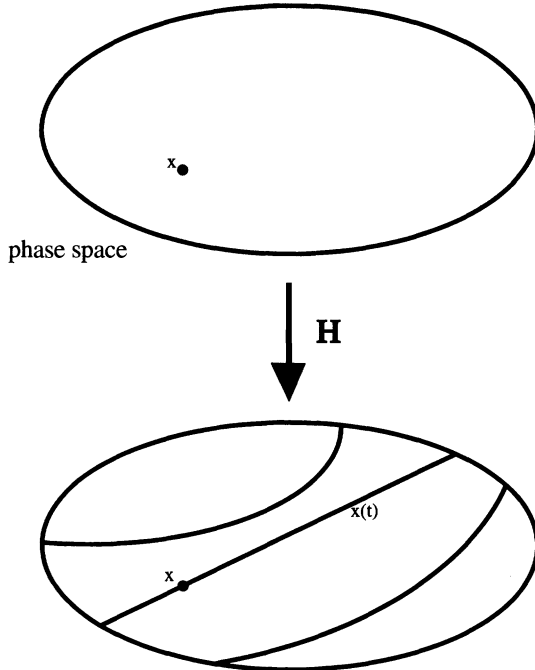


Fig. 1. Hamiltonian systems.

a representation of a Hamiltonian system: at the top, we have a phase space; specifying a Hamiltonian serves to determine a unique curve, $t \mapsto x(t)$, through each point x .

The simplicity of the Hamiltonian formalism makes the following *literal* approach to interpretation quite attractive. Given a Hamiltonian system, one would like to set up a bijection between points of phase space and dynamically possible states of the system. Then the theory at hand will be deterministic: given a point representing the present state of the system, the dynamical trajectory passing through that point represents the only physically possible past and future of the state.

Typically, it will be quite straightforward to develop such an interpretation. Most of the phase spaces of classical mechanics have the following form. One begins with a space Q , called the configuration space, which represents the possible configurations of some set of particles or fields relative to an inertial frame. One then constructs the cotangent bundle, T^*Q , of Q . This is the set of pairs (q, v) where $q \in Q$ and v is a vector at q . There is a canonical way of endowing T^*Q with a symplectic structure, so that it may be viewed as a phase space. Since a point $q \in Q$ represents a possible (generalized) position of the system, we can think of v as representing the system's (generalized) momentum.

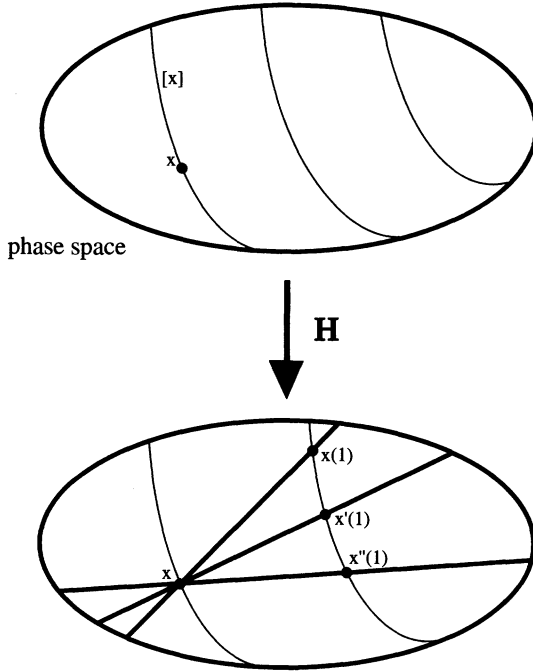


Fig. 2. Gauge systems.

The value of the Hamiltonian, H , at a given point of phase space will just be the energy of the state represented by that point.

(ii) *Gauge systems*

The geometry of the phase space of a gauge system is determined by a *presymplectic form*. This notion of geometry is weaker than the symplectic geometry of Hamiltonian systems.⁵ For our purposes, the upshot is the following: the phase space of a gauge system has a natural partition into subspaces, called *gauge orbits* (see the top half of Figure 2).⁶ The gauge orbits are all of the same dimensionality. Each point, x , of phase space lies in exactly one gauge orbit, denoted $[x]$. As in the Hamiltonian case, we specify the dynamics by choosing a real-valued function on phase space, the Hamiltonian.⁷ Whereas in the Hamiltonian case there was a single dynamical trajectory through each point

⁵ A presymplectic form, σ , on a manifold N is a closed but possibly degenerate two-form. It is standard to assume that the dimensionality of the null space of σ is the same at all points of N . A Hamiltonian system is a gauge system for which the null space is everywhere zero-dimensional.

⁶ Two points lie in the same gauge orbit iff they can be connected by a curve, all of whose tangent vectors are null vectors of σ . That is, the gauge orbits are constructed by integrating the null distribution of σ .

⁷ That is, we again look at the integral curves of vector fields which solve $X_H \lrcorner \omega = dH$. We require that H be 'gauge invariant' (see below).

of phase space, we find in the gauge-theoretic case that there are infinitely many trajectories through each point of phase space.⁸ The saving grace is that the dynamical trajectories through a given point, although disagreeing in general about which point represents the future state of the system at a given time, do agree about the gauge orbit in which this point lies. That is: if $t \mapsto x(t)$ and $t \mapsto x'(t)$ are dynamical trajectories which have their origin at the same point $x(0) = x'(0) = x_0$, then we have that $[x(t)] = [x'(t)]$ for all $t \in \mathbf{R}$, even when $x(t) \neq x'(t)$ (see Figure 2). Thus, although the presymplectic geometry is not strong enough to determine a unique dynamical trajectory through each point, it is strong enough to force all of the dynamical trajectories through a given point to agree about which gauge orbit the system occupies at each time.

There is an obvious obstacle to the application of gauge theories in classical mechanics. Hamiltonian systems have well posed initial-value problems: if we specify initial data (i.e. a point, x_0 , in phase space), then there is a unique solution of the equations of motion for that initial data (i.e. a dynamical trajectory, $x(t)$). Now, note that any given observable classical quantity can be represented by a real-valued function, f , on phase space (since the state of the system determines the values of all classical quantities). Thus, if we want to know the value of the quantity at time t_1 , we need only calculate $f(x(t_1))$, where $x(t)$ is the unique dynamical trajectory through x_0 . But in the case of a gauge system, there are many dynamical trajectories, $x(t)$, $x'(t)$, $x''(t)$, ..., through each point of phase space, so it is impossible to predict the future value of an arbitrary function on phase space from the initial state: in general, $x(t_1) \neq x'(t_1)$, so we expect that $f(x(t_1)) \neq f(x'(t_1))$.

But we know that it must be possible to apply such theories, since many of the most interesting classical field theories—electromagnetism, general relativity, Yang-Mills—are gauge theories.⁹ The solution is quite straightforward: observable quantities must be gauge-invariant. That is, if a real-valued function on phase space is to represent an observable quantity, then we require that f be constant on gauge orbits—if $[x] = [y]$, then $f(x) = f(y)$. The initial value problem of such a quantity is well-posed: if x_0 represents the initial state of the system, and $x(t)$ and $x'(t)$ are dynamical trajectories through x_0 , then $f(x(t_1)) = f(x'(t_1))$. Thus, despite the ambiguity in the evolution of the states of our gauge system, there is no ambiguity in the evolution of observable quantities—so long as we restrict our attention to gauge-invariant quantities.

⁸ Indeed, let X and X' be vector fields on N , and suppose that X solves $X \lrcorner \omega = dH$. Then X' solves $X' \lrcorner \omega = dH$ iff $Y = X - X'$ is a null vector field of σ (since $0 = X \lrcorner \sigma - X' \lrcorner \sigma = Y \lrcorner \sigma$).

⁹ There are also gauge theories which describe the gravitational interaction of point particles. These are of philosophical interest because they delimit the precise sense in which Mach's Principle can be implemented in classical physics. See Barbour and Bertotti [1982] and Lynden-Bell [1995].

But how can we interpret this formalism? There are three important interpretative approaches.

The most straightforward option is to mimic the *literal* approach that worked so well for Hamiltonian systems. This has the advantage of simplicity: one insists that each point of the phase space corresponds to exactly one physically possible state of the system. There is, however, a *prima facie* grave disadvantage to this approach: if the present state of the system is x_0 , then in general $x(t_1)$ and $x'(t_1)$ will correspond to distinct possible states of the system; the present will have many possible futures. Thus, literal interpretations render the theory indeterministic.¹⁰ Of course, if we supplement this account of the ontology of the theory with an account of measurement which implies that its observable quantities are gauge-invariant, then the indeterminism will not interfere with our ability to derive determinate predictions from the theory.

The second interpretative option is to stipulate that each gauge orbit of the phase space represents exactly one physically possible state. In this *simply gauge-invariant* case, our theory will be deterministic. We know that if $x(t)$ and $x'(t)$ are two dynamical trajectories for the same initial data, then $[x(t)] = [x'(t)]$ for all times t . Under a simply gauge-invariant interpretation, this is equivalent to saying that the present state determines the future and past states. In effect, our theory is construed as a theory of gauge orbits rather than points. There is an elegant construction which makes this precise: one can (usually) endow the set of gauge orbits of the gauge system with a symplectic structure and a Hamiltonian.¹¹ The resulting Hamiltonian system is called the reduced phase space. A point in the reduced phase space corresponds to a gauge orbit of the original gauge system; a dynamical trajectory of the reduced phase space tells us which gauge orbits the system passes through, given its initial gauge orbit; a real-valued function on the reduced phase space corresponds to a gauge-invariant function of the original system. Thus, the reduced phase space captures all of the gauge-invariant information of the gauge system. A simply gauge-invariant interpretation of the gauge system is equivalent to a literal interpretation of the reduced phase space—in both cases, one contends that physical possibilities stand in one-to-one correspondence with the gauge orbits of the original gauge system.

The third option is to adopt a *coarse-grained gauge-invariant* interpretation, according to which the representation relation between gauge orbits and physically possible states is many-to-one. Of course, since there is no physical difference between points in the same gauge orbit, coarse-grained gauge-invariant interpretations are deterministic (by the argument of the preceding paragraph).

¹⁰ The hole argument of Earman and Norton [1987] is a special case of this observation. See Belot and Earman [1998a, b].

¹¹ In rare cases, technical complications can make it impossible to carry out this construction.

This taxonomy provides us with a framework for discussing the interpretative possibilities for gauge theories. We can now see that it is an immediate consequence of the formalism of gauge theories that every such theory admits multiple interpretations. Settling on an interpretation is an important part of understanding the theory, since rival interpretations will disagree about whether or not the theory is deterministic, and about how much of the structure of the phase space is physically relevant.

This is likely to seem somewhat overwrought at this stage. After all, isn't it fairly clear that the preferred interpretation of a given gauge theory is a simply gauge-invariant interpretation? Since this is the same as a literal interpretation of the reduced phase space—and given that it is supposed to be straightforward to construct a literal interpretation of a Hamiltonian system—doesn't it follow that it should be straightforward to formulate such an interpretation? If this is so, then literal and coarse-grained gauge-invariant interpretations of gauge theories will never arise in practice.

I grant that simply gauge-invariant interpretations are, *ceteris paribus*, to be preferred. But I maintain that things are not always equal: sometimes it turns out that the available simply gauge-invariant interpretations are less plausible than their competitors. We will see an example of this below, in the case of electromagnetism. The problem is that, although the task of finding a simply gauge-invariant interpretation for a gauge theory does indeed reduce to the task of finding a literal interpretation of the theory's reduced phase space, it does not follow that it is straightforward to find a plausible interpretation of this Hamiltonian system. Literal interpretations of Hamiltonian systems are impeccable when the phase space has the structure of the set of possible positions and momenta of some set of fields or particles. But, in general, there is no guarantee that the reduced phase space of a gauge system will have such a structure.

3 Interpreting electromagnetism

The language of gauge theories provides the setting for a very elegant formulation of electromagnetism. Let physical space be modelled by some three-dimensional Riemannian geometry, S .¹² Then the phase space of electromagnetism consists of pairs, $(A(\xi), E(\xi))$, of vector fields on S , subject to the condition that $\text{div } E = 0$. That is, each point in the phase space of electromagnetism gives us a pair of maps, A and E , each of which assigns a three-vector to each point of S , with $\text{div } E = 0$. We call A the *vector potential*, and E the *electric field*. Our infinite dimensional phase space comes equipped with a natural presymplectic structure.

¹² The metric structure of S plays a hidden role in what follows: it allows us to define div , grad , and curl for non-Euclidean spaces.

The gauge orbits of the phase space have the following structure: (A, E) and (A', E') belong to the same gauge orbit iff for all $\xi \in S$, we have that $E(\xi) = E'(\xi)$ and that there exists a smooth function $\Lambda: S \rightarrow \mathbf{R}$ such that $A(\xi) = A'(\xi) + \text{grad } \Lambda(\xi)$. The Hamiltonian for electromagnetism is just $H = \int_s (|E|^2 + |\text{curl } A|^2) d\xi$. The dynamical trajectories are then determined by the following equations:

$$\left. \begin{aligned} \dot{A} &= -E \\ \dot{E} &= \text{curl}(\text{curl } A) \end{aligned} \right\} (3.1).$$

These are Maxwell's equations. Of course, these equations do not uniquely determine the evolution of the variables of electromagnetism: they do so only up to gauge. Thus, if we fix an initial point in phase space, (A_0, E_0) , then we find that the value of E is determined for all times, past and future. But the value of A is only fixed up to the addition of the gradient of a scalar function on space: if $(A(t), E(t))$ and $(A'(t), E'(t))$ are two solutions for our given initial data, then for all t , $E(t) = E'(t)$ and there exists a scalar $\Lambda(t)$ so that $A(t) = A'(t) + \text{grad } \Lambda(t)$. Maxwell's equations do not determine the future value of $A(t)$, but they do determine the gauge orbit in which $A(t)$ lies.

I will present three interpretations of electromagnetism, corresponding to the three strategies for interpreting gauge theories which were canvassed in the preceding section.¹³ It is helpful to have in mind some desiderata that we would like any interpretation of electromagnetism to fulfil.

First of all, of course, we would prefer that our interpretation render the theory deterministic: that is, we would like to find a gauge-invariant interpretation of electromagnetism.

Second, we would like our theory to be local. Here I draw a distinction between two types of locality:¹⁴

(I) *Synchronic locality*: the state of the system at a given time can be specified by specifying the states of the subsystems located in each region of space (which may be taken to be arbitrarily small).

(II) *Diachronic locality*: in order to predict what will happen *here* in a finite amount of time, Δt , we need only look at the present state of the world in finite neighbourhood of *here*, and the size of this neighbourhood shrinks to zero as $\Delta t \rightarrow 0$.

The first principle is meant to express the intuition that the properties of classical physical systems should be reducible to the properties of their

¹³ There are a few other interpretations of electromagnetism—see Brown [1994], Cao [1988], and Kennedy [1993]. None of them serves to undercut the themes developed below.

¹⁴ These are related to, but not identical with, the notions of separability and locality employed in the literature on the Bell inequality.

constituent parts. Classical fields, thought of as assignments of properties to the points of physical space, are paradigm examples of synchronically local systems—at each *here and now* the field has a state, and the state of the field itself is nothing but the sum of its states at each of these *here and nows*. Thus it makes sense, for instance, to speak of the state of the field in a given region, without reference to anything far removed from the region under consideration.

An object is local in the (strictly stronger) diachronic sense if its evolution in time is such that in order to know what its future state will be *here*, we need only know its present state in some finite *hereabouts*. Not every classical field is diachronically local. The Newtonian gravitational field is a well-known example of a diachronically nonlocal object—since gravitational effects propagate with infinite velocity, it is necessary to know the gravitational state *everywhere* in order to know exactly what will happen next *here*.

On the other hand, one expects electromagnetism to be both synchronically and diachronically local. It is, after all the theory of the electromagnetic *field* and so should be synchronically local. Furthermore, Maxwell's equations determine that electromagnetic radiation propagates at a fixed speed. This seems to imply immediately that electromagnetism is a diachronically local theory: since it will take some known finite amount of time for influences over there to reach here, I need only take into account what is happening over there when reckoning what will happen here, if I am interested in sufficiently large Δt .

We will see, however, that electromagnetism is diachronically local under only one of the three interpretations discussed below. Even worse, the theory is not even synchronically local according to one of these interpretations.

3.1 Interpretations

(1) *The vector potential as a physical field.* Under this first interpretation, one maintains that the vector potential, A , represents a physically real field on physical space. Most dramatically, one can maintain that the vector potential represents the velocity of a material ether.¹⁵ Then the electric field, $E = -\dot{A}$, would correspond to the acceleration of the material ether.¹⁶ This gives us a literal, hence indeterministic, interpretation of the gauge-theoretic formulation of electromagnetism: each pair (A, E) satisfying $\text{div } E = 0$ represents a distinct

¹⁵ Frank Arntzenius points out that this ether would be a strange object, since it would not possess the usual conserved quantities of classical fluid mechanics.

¹⁶ This is only one step removed from historical reality: in Maxwellian electrodynamics the current was sometimes interpreted as representing the acceleration of the ether. See Buchwald [1985], p. 24.

dynamical state of the ether, and two solutions, $A(t)$ and $A'(t) = A(t) + \text{grad } \Lambda(t)$, for the same initial data represent two physically distinct physical histories of the ether (according to A , *this* bit of ether ends up *here*; according to A' it ends up *over there*).¹⁷

This interpretation is synchronically local, since the state of the field reduces to the state of the field at each point of space. But it is diachronically nonlocal: changing the magnetic field in a given region can change the vector potential throughout space instantaneously—changes in the vector potential can propagate with infinite velocity even though electric and magnetic radiation propagates with a finite velocity. So in order to know what will happen to us next, we need to know what is presently happening arbitrarily far away. We will see an example of this phenomenon in the next section.

(2) *The traditional interpretation.* One would like to avoid the conclusion that electromagnetism is indeterministic. Thus, one would like to look for a gauge-invariant interpretation of electromagnetism. Here the most obvious option is the familiar one. We begin by defining the magnetic field to be $B \equiv \text{curl } A$. The value of the magnetic field at a point of physical space is a gauge-invariant quantity, since $\text{curl } A = \text{curl } (A + \text{grad } \Lambda)$. Together, E and B capture almost all of the gauge-invariant content of electromagnetism.¹⁸ Indeed, it is not difficult to show that the equations (3.1), together with the constraint $\text{div } E = 0$ and the identity $B \equiv \text{curl } A$, are equivalent to the familiar vacuum Maxwell equations:

$$\left. \begin{array}{ll} \dot{B} = -\text{curl } E & 2\text{div } B = 0 \\ \dot{E} = \text{curl } B & \text{div } E = 0 \end{array} \right\} \quad (3.2).$$

The initial-value problem for (3.2) is well posed: if we specify E and B at an initial time, then there is a unique solution of (3.2) which gives E and B for all future times.

Thus, if we stipulate that the ontology of electromagnetism consists of physically real electric and magnetic fields, then we have an interpretation which is deterministic (being gauge-invariant) and is clearly synchronically local (being based on fields).¹⁹ Furthermore, this interpretation is diachronically local: Maxwell's equations imply that the electric and magnetic fields propagate at the speed of light—in order to know what will happen here in Δt seconds, we need only consider the electromagnetic state in a sphere of radius $\Delta t \times c$.

¹⁷ As usual, we can none the less maintain that the theory is predictable if we can argue that only gauge-invariant quantities are measurable. If, for instance, our ether were an imponderable fluid, which interacted with ordinary matter only via electric and magnetic phenomena, then we would be unable to distinguish empirically between points of phase space which lie in the same gauge orbit.

¹⁸ The significance of this 'almost' will be made clear below.

¹⁹ Under this regime the vector potential becomes a useful mathematical fiction, with no physical content beyond that encoded in the magnetic field.

There remains one further question: is this interpretation simply gauge-invariant or coarse-grained gauge-invariant? It turns out that the answer depends upon the topology of physical space, S . If the topology of space is trivial, in the sense that S is simply connected, then the reduced phase space of electromagnetism is just the set of divergence-free electric and magnetic fields on S .²⁰ In this case, our interpretation may be viewed either as a literal interpretation of the reduced phase space or as a simply gauge-invariant interpretation of the original gauge system. If, however, S is topologically nontrivial, then the reduced phase space has a more complex structure (see below). Then the traditional interpretation counts as coarse-grained gauge-invariant: a number of gauge orbits correspond to each configuration of the electric and magnetic fields. In this case, there will be points in phase space, (A, E) and (A', E') , such that $[(A, E)] \neq [(A', E')]$ but A and A' correspond to the same magnetic field. We will see in the following section that this is the downfall of this otherwise very attractive interpretation.

(3) *Holonomies*.²¹ One would like an interpretation of electromagnetism which would be simply gauge-invariant, no matter what the topology of physical space. This is possible, but requires some ingenuity—and some sacrifices. The first step is to observe that although the value of the vector potential at a given point of physical space is not gauge-invariant (since in general $A(\xi) \neq A(\xi) + \text{grad } \Lambda(\xi)$), the integral of A around a closed curve, γ ,

$$h(\gamma) = \exp\left(\oint_{\gamma} iA_a(\xi)d\xi^a\right)$$

is gauge-invariant. In fact, if we substitute a vector potential, A' , into the integral which defines $h(\gamma)$, then this quantity is unchanged *if and only if* A' is in the same gauge orbit as A . $h(\gamma)$ is called the *holonomy* around γ , and is a complex number of unit modulus.

It is a remarkable fact that we can construct the reduced phase space of electromagnetism by taking the set of holonomies around all curves in physical space as the elements of our configuration space. That is, if we call the set of closed curves in physical space *loop space*, each point in our configuration space is just a map from loop space to the complex numbers of unit modulus. Let us call these maps holonomy maps. The value of a given holonomy map on a given loop is just the holonomy around that loop. We proceed to construct the phase space by building the cotangent bundle, which is just the set of pairs consisting of a holonomy map and a divergence-free electric field. After we

²⁰ A space is said to be *simply connected* if every closed curve in the space may be contracted to a point without leaving the space. Otherwise, it is said to be *multiply connected*. Thus, in two dimensions the sphere and the plane are simply connected, while the cylinder and the torus are multiply connected (imagine drawing a circle which goes around the circumference of the cylinder or torus—such a curve cannot be contracted to a point without leaving the surface).

²¹ See Wu and Yang [1975] for an influential presentation of this approach.

impose the canonical symplectic structure, and the correct Hamiltonian, we end up with a Hamiltonian system which is the reduced-phase-space formulation of electromagnetism.

We would like to give a literal interpretation of this reduced phase space formalism. We know that this can be done in terms of the electric and magnetic fields iff physical space is topologically trivial. But in general, the topology of S is nontrivial, and this move is not available—we cannot identify the space of holonomy maps with a space of tensors on S .

It is, none the less, possible to formulate an interpretation of electromagnetism which is simply gauge-invariant no matter what the topology of space. The reduced phase space formulation of electromagnetism suggests that a state of the electromagnetic field should be thought of an assignment of complex numbers to closed curves in space (holonomy), together with an assignment of vectors to points of space (electric field). This requires a revision of our notion of field. We can no longer think of the electromagnetic field as simply being an assignment of properties to points of space. Rather, we must also consider closed curves in space to be carriers of the electromagnetic predicates. This interpretation is, of course, deterministic, since it is simply gauge-invariant. But it is also synchronically (and hence also diachronically) nonlocal, since specifying the electromagnetic state of any given region of physical space requires knowledge of the holonomy around every loop in space, and hence requires mentioning regions of space arbitrarily far away from the one under consideration.

These, then, are our three interpretations. Here we find ourselves in a situation of the sort alluded to at the end of the previous section. *Ceteris paribus*, we would prefer a simply gauge-invariant interpretation of our theory. But in the case at hand this is awkward: unless physical space has a very special topological structure, the reduced phase space of electromagnetism cannot be viewed as the space of positions and momenta of a set of fields relative to physical space. Thus, simple gauge-invariance can only be had at the price of a revision our intuitions about classical fields.

In this context, it seems clear that interpretation (2) is to be preferred. It is the only interpretation which is both deterministic and diachronically local. When S has a non-trivial topology, it is true, this traditional interpretation has a vice: since it is coarse-grained, it loses information which is stored in the other interpretations. As a result, this interpretation is vulnerable to empirical refutation: it is an empirical question whether or not there exist measurable physical quantities which distinguish between each pair of gauge orbits. If there were any such quantities, no coarse-grained gauge-invariant interpretation would be tenable. Within the realm of classical physics, however, (2) is vindicated—there are no phenomena which allow one to distinguish between two gauge orbits $[A]$ and $[A']$ which correspond to the same magnetic field.

Thus, there are no grounds internal to electromagnetism upon which to criticize the traditional interpretation.

4 Quantization and the Aharonov–Bohm effect

The story of the Aharonov–Bohm effect is normally told in a very abbreviated form: something like: ‘The quantum treatment of a charged particle in an external magnetic field shows that it is possible to distinguish between vector potentials which correspond to the same magnetic field. So quantum mechanics shows that the vector potentials of electromagnetism are physically real.’ This sort of account gives the correct flavour. But it is also quite misleading: if the vector potential is physically real in the way that the electric field is supposed to be, then electromagnetism is indeterministic (as in interpretation (1) of the previous section)—but few commentators mean to commit themselves to such a view. In this section, I present a rational reconstruction of the way in which the Aharonov–Bohm effect bears upon the considerations of the previous section.²²

The story has a complex structure, so I have broken it down into several parts. The first part consists of a generic account of how interpretative beliefs about classical theories bear upon quantization, and vice versa. In the second, this apparatus is applied to the special case of a charged particle in a static magnetic field. The Aharonov–Bohm effect makes its appearance in the third and final part.

4.1 Quantization and interpretation

Here I focus on *canonical quantization*, where the initial input is a classical theory in Hamiltonian form, and the output is a quantum theory which—one hopes—has the original classical theory as its $\hbar \rightarrow 0$ limit. It is important to emphasize that neither quantization nor the taking of a classical limit is entirely straightforward. The details of either sort of process vary greatly from case to case. Thus there is considerable play at either end whenever we attempt to set up a correspondence between classical and quantum systems. This freedom in bridging the classical–quantum divide is closely related to an interplay which exists between interpretative issues at the classical and quantum levels.

It is well known that all infinite dimensional Hamiltonian systems have infinitely many quantizations. That is: there are infinitely many unitarily inequivalent quantum field theories which are quantizations of any given classical field theory. This fact has received some attention from philosophers

²² See Healey [1997] for a complementary discussion of the bearing of the Aharonov–Bohm effect upon approaches to interpreting quantum mechanics.

of physics. But it has tended to be regarded as one pathology among many in the foundations of quantum field theory. What has received little (no?) attention from philosophers, is the fact that the same situation is endemic among finite dimensional systems. Indeed, any Hamiltonian system with a topologically nontrivial phase space has infinitely many quantizations.²³ Thus, the ambiguity involved in quantization cannot be dismissed as a pathological feature of a theory which is still under development—it is to be found everywhere, even in ordinary, non-relativistic quantum mechanics.

We have to ask ourselves what attitude to take towards inequivalent quantizations of a given classical theory—how many of them are physically significant? Often there will be a distinguished quantization.²⁴ In this case it will be tempting to maintain that it alone should count as an admissible quantum treatment of the phenomenon under investigation. On the other hand, it can happen that one has reason to believe that there are many acceptable quantizations of the system (each appropriate for modelling a distinct physical situation), or that the structure of the classical system does not single out a preferred quantization (as in quantum field theory on curved spacetime).

Ultimately, of course, one hopes that empirical considerations will determine which quantizations should be taken seriously. But one sometimes finds oneself in a situation where available data underdetermine the question—indeed this situation is probably the norm at the frontiers of theoretical physics. Here there are a number of epistemic resources which could be mobilized to fill the gap. Among these are interpretative beliefs about other theories.

One's interpretative beliefs can shape one's judgements as to the relevance of certain quantizations or approaches to quantization. Conversely, one must accept that one's interpretative beliefs are open to revision in light of the empirical success of the approaches to quantization which they suggest. If one's interpretation of a given classical theory suggests that quantization A is superior to quantization B, then one is bound to revise one's interpretative judgements if it turns out that A is empirically untenable.

This is, I maintain, the best way of thinking of the import of the Aharonov–Bohm effect: we come to quantum mechanics with prior interpretative beliefs about electromagnetism; these suggest a particular approach to quantizing classical systems involving electromagnetic fields; when this does not pan out, we are obliged to revise our understanding of electromagnetism.

²³ See Woodhouse [1980] for details. The physically inequivalent quantizations of a given Hamiltonian system are parameterized by the cohomology group $H^1(M, U(1))$, where M is the classical phase space, and $U(1)$ is the group of complex numbers of unit modulus. This group is a topological invariant—it is determined by the topology of the classical phase space. For our purposes, it suffices to assume that this cohomology group is nontrivial iff the classical phase space is multiply connected, although this is not quite true.

²⁴ Indeed, this is always the case for quantizations of finite dimensional systems, since the group $H^1(M, U(1))$ always has a preferred element—namely, the identity.

4.2 Quantizing the charged particle in a magnetic field

The Aharonov–Bohm effect involves a quantum charged particle moving in a static magnetic field. So we are interested in quantizing the classical treatment of such a particle.²⁵

In order to construct the classical treatment of the particle, we proceed as follows. We let S be the region of physical space in which the particle is free to move. The configuration space of the particle is isomorphic to S . The phase space of the particle is just the set of possible positions and momenta for such a particle (the cotangent bundle of S), endowed with a symplectic structure which differs from the symplectic structure for a free particle by a factor which is determined by the magnetic field. Finally, the Hamiltonian, as in the treatment of a free particle, is just the kinetic energy.

This system will have a unique quantization iff S is topologically trivial. We will be interested in a case where S is (homeomorphic to) ordinary Euclidean three-dimensional space with the y -axis removed, and the electric and magnetic fields are zero. Thus, our model of the charged particle moving in S will be identical to our model of a free particle moving in S . This model has infinitely many quantizations. These quantizations can be put in a one-to-one correspondence with the set of complex numbers of unit modulus.²⁶ That is, specifying a complex number, z , with $|z| = 1$ determines a quantization $Q(z)$ of our classical system. Of course, one possibility stands out: the quantum system $Q(1)$. Intuitively, the quantization $Q(z)$ predicts that a quantum charged particle which is transported around the y -axis will have its phase shifted by a factor of z .²⁷

Now, in the absence of empirical data, what attitude should one take towards these various quantizations? If we believe in a simply gauge-invariant or literal interpretation of electromagnetism, we will be open to the idea that each of the quantizations of the charged particle represents a genuine physical possibility. After all, the Hamiltonian representation of this particle contains only partial information about the electromagnetic state: the Hamiltonian H contains no information at all, and the symplectic form contains information about the magnetic field alone. But in the case at hand, where S is multiply connected, we know that specifying B is inadequate—the true electromagnetic state of the world contains considerably more information, since many gauge orbits of the phase space correspond to the same B . In the present case this excess information can be encoded as a single complex number of unit modulus—the holonomy around a loop which circumnavigates the y -axis.²⁸ In light of this fact, it would be quite natural for someone who

²⁵ I.e. the particle is treated quantum-mechanically while the field is treated classically.

²⁶ Since our phase space is T^*S and $H^1(T^*S, U(1)) = U(1)$.

²⁷ We would have to be a bit more careful at this point if we were working with $B \neq 0$, and specify a particular curve which circumnavigates the axis. But the points established below would still go through.

espouses interpretations (1) or (3) of electromagnetism to adopt the following line. In order to specify the electromagnetic state of the world, one must specify the holonomy around a loop around the y -axis as well as the values of E and B . Since our Hamiltonian treatment of the charged particle in a static magnetic field does not contain this additional information, it is no surprise that there are multiple quantizations and that these are in one-to-one correspondence with the possible values of the holonomy: this just means that there is one quantization for each possible electromagnetic state with $E = B = 0$. Each of these quantizations should be taken seriously, since each is appropriate for a distinct physical situation.

If, on the other hand, we accept the traditional interpretation of the theory—according to which the electric and magnetic field are the only physical realities, and any information contained in the reduced phase space of electromagnetism which is not determined by the values of these fields is descriptive fluff—then we will not believe that there are any hidden variables which could determine which of the quantizations of our classical system is the physically correct one. Our attitude will be: having built all of the physically relevant information into our classical model of the particle moving in a magnetic field, we have no reason to believe that it admits of more than one physically realistic quantization. This line of argument will seem especially convincing when it is recalled that we use the same classical model to represent an uncharged free particle in region S —and of course we believe that in this latter case there is a unique physically correct quantization.

Thus the coarse-grained gauge-invariant interpretation (2) leads to a different approach to quantizing the charged particle in a magnetic field than that suggested by the other interpretations of electromagnetism. And here we have an empirical question: what would a charged quantum particle do in a situation like this? If we find that distinct gauge orbits of vector potentials which correspond to different holonomies but to the same magnetic field lead to different interference patterns, then the coarse-grained gauge-invariant interpretation is sunk—it leads to an empirically inadequate quantum theory.

4.3 The Aharonov–Bohm effect

It is not possible to subject interpretation (2) to a direct empirical test—we simply do not know enough about the topology of physical space. And even if

²⁸ This is because when $E = B = 0$, all loops which are homotopic (i.e. which can be continuously deformed into one another) have the same holonomy. And when $S = \mathbb{R}^3 \setminus \{y\text{-axis}\}$, the equivalence classes of homotopic loops are parameterized by the number of times that they loop around the y -axis. In our case, this means that specifying the holonomy of any curve which loops around the y -axis once suffices to determine the holonomy of all curves, since the holonomy of a curve which wraps around the y -axis n times is just n times the holonomy of a curve which wraps around the axis once.

we did, it is unlikely that we would be able to detect these sorts of effects on the cosmological scale. There is, however, a tabletop experiment which is widely regarded as refuting interpretation (2). This experiment was first suggested by Aharonov and Bohm in 1959, and was carried out shortly thereafter. The idea is to restrict an electron to a region which can be treated for all practical purposes as being topologically nontrivial. One should then be able to see which quantizations are physically relevant.

The experiment works as follows. One constructs a solenoid—a conducting wire coiled around a cylinder—whose length is long compared to the wavelength of the particle under consideration. When current runs through the solenoid, a magnetic field is created inside the device, but the external magnetic field is unaffected.²⁹ Now suppose that we want to model the behaviour of a classical charged particle in the field-free region external to, but near, a solenoid. What manifold should we use as our configuration space? We could use S' , some region of physical space which includes the solenoid. But then our symplectic form will have to encode information about the state of the magnetic field inside the solenoid. This will be formidably complicated when the solenoid is operational.

There is, however, another option for describing the system which promises to be much simpler. We can assume that the walls of the solenoid are impenetrable. Now, under any reasonable interpretation of electromagnetism, the magnetic field is synchronically local (more on this below). So it seems that we should not have to know anything about the interior of the solenoid in order to model the behaviour of the particle. Thus, we should be able to take S as our configuration space, where S is just S' with the points occupied by the solenoid deleted. For all practical purposes, we can treat the solenoid as being infinitely long. The resulting phase space is of the form discussed in subsection (2) above. Since the classical particle is indifferent to the state of the field inside the solenoid, we always use the same Hamiltonian system, no matter how much current is running through the solenoid.

When we want to construct our quantum theory of the charged particle moving in the region external to the solenoid, we find ourselves in the situation described in the previous subsection: the inequivalent quantizations of our system are parameterized by the complex numbers of unit modulus, and we want to know how many of them model physical possibilities. But now we have a manageable experimental question: what happens when we shoot a beam of electrons by a solenoid? Experiment reveals that the interference pattern which results *does* depend on whether or not the solenoid is operational. In fact, it

²⁹ It follows that the vector potential propagates with infinite velocity: the holonomy around a closed curve which loops around the solenoid once is equal to the magnetic flux through the solenoid; if we switch the thing on, then the values of the vector potential at some point arbitrarily far away must change instantaneously.

turns out the quantization which yields the correct predictions for a given experimental situation is $Q(z)$, where z is the holonomy of a loop which circumnavigates the solenoid.

Thus interpretation (2) of electromagnetism is in trouble. It led us to expect that only one of the quantizations would be empirically interesting. This turns out to be false. This, in turn, casts considerable doubt on the position that the electric and magnetic fields encode all of the physically relevant information. Interpretations (1) and (3), on the other hand, allow that the holonomies around closed curves carry genuine physical information over and above that contained in a specification of the fields. They are thus able to easily account for the Aharonov–Bohm effect.

Where does this leave interpretation (2)? Above, I noted that it is permissible to employ S as the configuration space of the classical system only if we believe that the behaviour of the classical charged particle depends only on the field values where it is, and not on the field values in the interior of the solenoid. Now, it is possible to alter interpretation (2) so that the magnetic field acts where it is not. In doing so, one gives up synchronic—and hence also diachronic—locality. But one is able to thus account for the Aharonov–Bohm effect. Presumably, few would be happy with this nonlocal version of interpretation (2).

For the vast majority of physicists, this has provided sufficient reason to turn away from interpretation (2), and to take seriously Aharonov and Bohm’s conclusion that we should regard the vector potential ‘as a physical variable’ ([1959], p. 491). But what exactly does this mean? Is it an endorsement of interpretation (1) or of interpretation (3)? Aharonov and Bohm themselves seem to have had interpretation (1) in mind.³⁰ This was, however, before it was shown that one could give a simply gauge-invariant interpretation in terms of holonomies. Today I think that almost anyone who takes the care to distinguish between interpretations (1) and (3) would endorse the latter: people seem to be more willing to make the move to thinking of fields as properties of loops—thus sacrificing even synchronic locality—than they are to sacrifice determinism.³¹

5 Conclusion

The articulation of the content of our physical knowledge is one of the chief tasks of philosophy of physics. Much of this work is interpretative in nature.

³⁰ They gloss their comment quoted above by saying that ‘[t]his means that we must be able to define physical difference between two quantum states which differ only by a gauge transformation.’ I read this as indicating that they believe that points in the same gauge orbit correspond to distinct physical situations.

³¹ One factor here is surely that holonomies and loops have provided a fruitful framework for the quantization of gauge theories such as electromagnetism, general relativity, and Yang–Mills theories. This sits well with the theme developed in the next section.

Naïvely, one seeks the correct interpretation of a given physical theory X. But, of course, all extant physical theories are *false*—none is both fully relativistic and fully quantum. What sense does it make to speak of the *correct* interpretation of a false theory?

A better picture would be this. We can begin by thinking of the content of a physical theory as consisting of the set of worlds at which it is true. This provides only a first approximation to the notion of physical content: we certainly do not view each interpretation as being on a par. So we should think of the content of a physical theory as the set of worlds described by the theory together with some additional structure which encodes our evaluation of the relative merits of each of the possible interpretations. Interpretation is the articulation of this further structure. Sometimes, as in electromagnetism, there will be a single most preferred interpretation. We may, if we like, think of this interpretation as being the ‘correct’ one. We will, of course, reserve the right to revise such judgements.

The formalism of the theory picks out the set of possible worlds which underlies the content of the theory. What determines the further structure? Here purely metaphysical views will play some role. But surely our beliefs about the structure of our own world make a large contribution here. We feel that there is a sense in which a world containing a physically real ether is less like our own than one in which fields are free-standing. This contributes to our sense that some interpretations of electromagnetism are more far-fetched than others.³² To the extent that such interpretative judgements place constraints on our beliefs about where the actual world might sit in the space of possible worlds, they are indeed judgements about our world. There is a clear sense, then, in which the interpretation of false theories teaches us about this world. Our beliefs about our world are reflected in our understanding of our false physical theories; so getting clear on the content of a false theory is one way to make explicit our beliefs about our world. Admittedly, this is a strange way to learn about the world. But it is also a fruitful one for us: in the absence of a true theory, our false theories provide much of our understanding of the structure of the world.

The discussion of the preceding sections shows that distinct interpretations of electromagnetism constitute different ways of understanding the theory: according to some interpretations the theory is deterministic, according to others it is nonlocal. The significance of these divergent approaches becomes clear when we examine the conceptual relationships between electromagnetism and quantum mechanics: different interpretations of the classical theory of the

³² This sort of local notion of verisimilitude does not require commitment to a corresponding global notion. I can understand what it means to say that an ether world is a remote possibility, without necessarily understanding what it would mean to ask: ‘Which is closer to the truth, Maxwell’s theory or quantum mechanics?’

electromagnetic field suggest distinct approaches to resolving the ambiguity inherent in the quantization of classical models of the behaviour of charged particles. Fairness then requires us to recognize that the empirical success of one or another of these approaches to quantization bears upon our understanding of electromagnetism. Similarly, the question of whether the general covariance of general relativity should be understood as a principle of gauge-invariance forms a bridge between the interpretative problems of general relativity (What is the nature of the existence of spacetime points?) and those of quantum gravity (Do time and change exist at the fundamental level?).³³ Different approaches to understanding the general covariance of general relativity are associated with different solutions to the interpretative problems of classical and quantum gravity. Thus, the empirical vindication of a given approach to quantizing gravity may have repercussions for our understanding of classical spacetime ontology.

Examples such as these force us to conclude that theories cannot be interpreted in isolation from one another: understanding intertheoretic relations is a crucial component of the articulation of the content of individual physical theories. In order to proceed, we need to develop a way of thinking about the structure and content of physics which provides a useful framework for discussion of intertheoretic issues, as well as intratheoretic ones. This should start from a recognition that we have a multitude of theories on the books: classical mechanics, statistical mechanics, electromagnetism, quantum mechanics, quantum statistical mechanics, quantum field theory, special relativity, general relativity, and so on. These theories tell us about very different worlds. Some are populated by particles, some by fields. In some the spatiotemporal structure is an unchanging backdrop, in some it is an active and changing participant. It is an important fact that we find that peaceful coexistence rather than competition is the rule, despite the nontrivial overlap between the domains of applicability of these apparently incompatible theories. Each of these theories informs us about our world, despite their profound divergence of opinion concerning ontology.

So we have a network of theories. This network is often described as a hierarchy, the idea being that some theories are more fundamental than others. In fact, a web or a lattice would be a more appropriate metaphor here, since theories often have more than one limit—special relativity is the $c \rightarrow 0$ limit of GR, while the Newton–Cartan theory is its $c \rightarrow \infty$ limit. It is possible to speak of one theory being more fundamental than another only so long as we don't make the mistake of assuming that 'more fundamental' gives us a linear ordering of the class of theories.

Each of our theories is empirically adequate within its own domain of applicability, but shares parts of this domain with other theories which save

³³ See Belot and Earman [1998a, b].

their own phenomena. Given this situation, it seems essential to demand for every pair of overlapping theories an assurance that their empirical predictions mesh in the appropriate manner. The correspondence principle, for example, can be understood along these lines: as requiring that quantum mechanics be able to account for the empirical adequacy of classical mechanics.³⁴

The following picture emerges. Each of our physical theories is part of a network of theories which stand in subtle relations to one another. Each theory contributes to our understanding of the world not only in virtue of its internal structure and empirical adequacy, but also because of the relations in which it stands to other theories. Indeed, we have seen that the content of a given theory may depend upon how it is situated in the network—quantum considerations help to fix the meaning of the terms of electromagnetism, the substantial–relational debate about the spacetime of general relativity is intimately connected with ongoing work on quantum gravity. In terms the web metaphor: if one looks only at individual nodes of the web, one gains only partial information; in order to grasp the full content of the fact that this web is useful for describing our world, one must also look at the way that nodes are situated with respect to their neighbours, and at the strands which join these nodes. More graphically—if somewhat distastefully—if one wants to survey the shape of an object caught in a net, it is useful to note how the threads bulge, as well as where they meet.

In the case of electromagnetism, the concrete payoff of this approach is that one can clarify the interpretative status of the theory by looking at the relation between the Hamiltonian system which models a classical charged particle in a magnetic field and its quantization. Constructing a new node in the web of theories (quantum mechanics) and linking this node to its neighbour (classical mechanics, via quantization and classical limits) has necessitated adjusting the position of the web elsewhere: we have been forced to relinquish an attractive interpretation of electromagnetism, in light of the conceptual structure of the relation between classical and quantum mechanics, and the results of empirical investigations.

How can we understand this constraint on interpretation? I propose that we think of it as a requirement that our preferred interpretations be *fruitful*. When we discover that holonomies are preferable to magnetic fields for interpreting electromagnetism, we do not, *strictly speaking*, discover something about the ontology of our world—in the strictest sense, it contains neither holonomies

³⁴ It is important to emphasize that one cannot expect too much of the correspondence principle and its generalizations. In general, as one takes a classical or non-relativistic limit the ontology of a given theory does *not* go over into the ontology of the limiting theory (Rohrlich [1988]). And this should not bother us if we keep in mind that both the given theory and the limit theory are false. Similarly, it is important to keep in mind that the correspondence principle requires only the claim that each of the theories is empirically adequate the consistency of, and not that the theories be empirically equivalent.

nor fields. But we do discover that one interpretation provides a more fruitful way of thinking about a certain range of phenomena than the other—both provide a satisfying picture of electromagnetism, but only one of them provides helpful hints about the quantum realm. For this reason, we think that electromagnetism interpreted in terms of holonomies is closer to the truth about our world than electromagnetism interpreted in terms of magnetic fields. This changes our opinion about the situation of the actual world in the space of possible worlds. It is in this sense that the interpretative fallout from the Aharonov–Bohm effect teaches us something about our world.

Acknowledgements

Earlier versions of this paper were presented at the University of Pittsburgh, Carnegie Mellon University, the University of Western Ontario, the University of Michigan, Stanford, and UCLA. I would like to thank Frank Arntzenius, Richard Boyd, Joe Camp, John Earman, Kit Fine, Fritz Rohrlich, Laura Ruetsche, and Lisa Shapiro for helpful conversations. I would also like to acknowledge the generous support of SSHRC, and the hospitality of the Center for Philosophy of Science at the University of Pittsburgh.

*Department of Philosophy
Princeton University
Princeton, NJ 08544
USA*

References

- Aharonov, Y. and Bohm, D. [1959]: ‘Significance of Electromagnetic Potentials in Quantum Theory’, *Physical Review*, 115, pp. 485–91.
- Barbour, J. and Bertotti, B. [1982]: ‘Mach’s Principle and the Structure of Dynamical Theories’, *Proceedings of the Royal Society of London, Series A*, 382, pp. 295–306.
- Belot, G. and Earman, J. [1998a]: ‘From Metaphysics to Physics’, forthcoming in H. Brown, J. Butterfield, and C. Pagonis (eds), *From Physics to Philosophy*.
- Belot, G. and Earman, J. [1998b]: ‘Pre-Socratic Quantum Gravity’, forthcoming in C. Callender and N. Huggett (eds), *Physics Meets Philosophy at the Planck Scale*.
- Brown, J. R. [1994]: *Smoke and Mirrors: How Science Reflects Reality*, New York, Routledge.
- Buchwald, J. Z. [1985]: *From Maxwell to Microphysics: Aspects of Electromagnetic Theory in the Last Quarter of the Nineteenth Century*, Chicago, University of Chicago Press.
- Cao, T. Y. [1988]: ‘Gauge Theory and the Geometrization of Fundamental Physics’, in H. R. Brown and R. Harré (eds), *Philosophical Foundations of Quantum Field Theory*, Oxford, Oxford University Press.

- Cohen, R. and Stachel, J. (eds) [1979]: *Selected Papers of Léon Rosenfeld*, Dordrecht, Reidel.
- Earman, J. and Norton, J. [1987]: 'What Price Spacetime Substantivalism? The Hole Story', *British Journal for the Philosophy of Science*, 38, pp. 515–25.
- Healey, R. [1997]: 'Nonlocality and the Aharonov-Bohm Effect', *Philosophy of Science*, 64, pp. 18–41.
- Kennedy, J. B. [1993]: 'Interpreting Gauge Symmetries: a Proposed Solution to the Problem of Gauge Invariance', unpublished.
- Lynden-Bell, D. [1995]: 'A Relative Newtonian Mechanics', in J. Barbour and H. Pfister (eds), *Mach's Principle: From Newton's Bucket to Quantum Gravity*, Boston, Birkhäuser.
- Rohrlich, F. [1988]: 'Pluralistic Ontology and Theory Reduction in the Physical Sciences', *British Journal for the Philosophy of Science*, 39, pp. 295–312.
- Woodhouse, N. [1980]: *Geometric Quantization*, Oxford, Oxford University Press.
- Wu, T. T. and Yang, C. N. [1975]: 'Concept of Nonintegrable Phase Factors and Global Formulation of Gauge Fields', *Physical Review D*, 12, pp. 3845–57.