

Synopsis and Discussion

Philosophy of Gauge Theory

Saturday-Sunday, 18-19 April 2009

Center for Philosophy of Science

817 Cathedral of Learning

University of Pittsburgh

Pittsburgh, PA USA

“Gauge Fields: What Isn't There?” Richard Healey, Philosophy, University of Arizona

“The Aharonov-Bohm Effect: Solved” James Mattingly, Philosophy, Georgetown University

“Intrinsic Geometrical Structure and Fiber Bundles” Tim Maudlin, Philosophy, Rutgers
Comments: Gordon Belot, Philosophy, University of Michigan

“Spontaneous Symmetry Breaking and the Higgs Mechanism “ Ward Struyve, Perimeter
Institute for Theoretical Physics

Discussants: Antigone Nounou, Philosophy, University of Minnesota, and Oliver Pooley,
Philosophy, Oxford University.

Contents

Contributions by

Antigone Nounou
Version 1. April 29, 2009.

John Earman
Tim Maudlin
Richard Healey
Version 2. May 2, 2009.

James Mattingly
Version 3, May 6, 2009

Tim Maudlin
Version 4, May 7, 2009

Richard Healey
Ward Struyve
Gordon Belot
Version 5, May 8, 2009

Tim Maudlin
Version 6, May 10, 2009

Ward Struyve/Antigone Nounou (revisions)
Version 7, May 12, 2009

Antigone Nounou
Version 8, May 15, 2009

James Mattingly
Version 9, May 20, 2009

Comments by Antigone Nounou

Richard: Loops, Holes and Causes

Richard Healey's talk was divided in two parts. In the first part he argued that we are not justified in believing that localized gauge potential properties are there, but we are in believing that holonomy properties are. In the second part, he conceded that the holonomy interpretation offers an incomplete local and causal account, but he maintained that the onus is on QM.

PART I

Despite appearances, the argument against gauge potential properties, which is also an argument for non-separable holonomy properties, is not analogous to the hole argument in GTR. Similarities aside (both arguments appeal to unobservability and both deny the existence of certain theoretical structures), the two arguments differ in significant respects: the hole argument is a *reductio ad indeterminism* that denies the existence of space-time points, whereas the argument against localized gauge potential properties consists in criticizing a bad abductive argument whose conclusion, that there exist gauge potential properties which are distributed over space-time points and are represented by the connection ω , is unsupported.

Since Yang-Mills theories contain two theoretical structures that may be thought to represent gauge potential properties, namely holonomies and connections whose local representatives are the gauge potentials, two possible interpretations are available: the holonomy interpretation (HI) and the localized property interpretation (LI). According to HI, it is the holonomies of ω that represent gauge potential properties. Thus, holonomy-equivalent connections, e.g. ω and ω' , represent the same distribution of gauge potential properties that attach non-separably on entire loops of space-time. LI, on the other hand, takes ω to represent gauge potential properties directly, and this implies that the holonomy-equivalent connections ω and ω' may represent different distributions of gauge potential properties over space-time.

The empirical facts from classical Yang-Mills theories, however, favor HI and not LI. This constitutes evidence that we should believe in the existence of holonomy properties but not of localized gauge potential properties. Moreover, the lack of empirical evidence for LI implicates empirical underdetermination which renders certain of its statements not merely unverifiable but also meaningless. For example, if ω and ω' are holonomy-equivalent connections, the statements " ω represents the distribution of localized gauge potential properties" and " ω' represents the distribution of localized gauge potential properties" are meaningless in the sense that they are not empirically incompatible statements. And even meaningful (but unverifiable) LI statements, like for example "there are localized gauge potential properties distributed over space-time points in some way compatible with the actual holonomy properties on space-time loops", cannot acquire meaning through the aforementioned meaningless statements, and cannot be justified by inference to the best explanation for two reasons. First, LI fails to account for the relations between holonomy properties; and second, an alternative explanation of these relations is available through loop supervenience.

PART II

The second part of the talk begun by positing two requirements that should be met in understanding the Aharonov-Bohm effect: non-local action should be avoided and an explanation that is causal or quasi-causal should be provided. The conclusion is that HI may fulfill the first but not the second, and the reason for this failing is to be found not in HI itself but in quantum mechanics which does not yield causal or quasi-causal explanations.

The propagation of non-separable holonomy properties conforms to relativistic locality. HI also conforms to local action provided that the latter only applies to space-like or null separated loops. On the other hand, holonomy properties act by affecting physical properties of matter, and they would act non-separably yet locally on classical charged matter fields, if such fields existed. The matter fields involved in the Aharonov-Bohm effect, however, are described by relativistic quantum mechanics, and it is a fact that quantum mechanical accounts are not causal. Therefore, the failure of HI to provide a causal account, which becomes evident in both the electric and the time-dependent magnetic flux Aharonov-Bohm effect, should come as no surprise. But this is not a weakness of HI, nor is it detrimental. For, physical theories need not and do not in general reveal the causes of phenomena.

Antigone's concern:

The modification of local action (so that it applies to space-like or null separated loops only) may turn out to be ad hoc. But for the time being she's got nothing more to say on that.

James: The A-B Effect: Solved

James Mattingly argued that there is an acceptable interpretation of electromagnetism according to which both the gauge potentials and electromagnetic bulk fields are mere calculational devices, but the basic quantities are the components of the electromagnetic fields that are generated by individual elements of charge-current. Since the evidence for this interpretation is the same with the evidence for classical electromagnetism, it raises no epistemic concerns; and it is metaphysically parsimonious because it only postulates a single entity that acts locally even when quantum particles are concerned.

The evidence for the field components produced by the charge-current elements comes from the Aharonov-Bohm effect, which indicates that the bulk field is zero in all regions that are accessible to quantum particles. Gauge potentials are non-zero somewhere in that region, but being ill-defined and non/deterministic, or non-local and non-separable, they do not constitute good ontological choices. A good gauge ontology would rather be well-defined, local, separable and ontologically robust. The Liénard-Wiechert (LW) potentials, which are uniquely specifiable without gauge fixing, make good candidates for such ontology and so do the component fields, the electric and magnetic fields that are generated by individual elements of charge-current.

A Lagrangian that employs the component fields and is appropriate for the Aharonov-Bohm effect is the Darwin Lagrangian for a charge a in the field produced by a collection of charges, those circulating in the solenoid. The last term in this Lagrangian,

$\frac{e_a}{2c^2} \sum_{b \neq a} \frac{e_b}{R_{ab}} [\mathbf{v}_a \cdot \mathbf{v}_b + (\mathbf{v}_a \cdot \mathbf{n}_{ab})(\mathbf{v}_b \cdot \mathbf{n}_{ab})]$, corresponds to the magnetic potential energy, and

despite the fact that it does not contain the magnetic field itself explicitly, some algebraic manipulation does the trick: the term turns out to be equal to

$$\frac{e_a}{2c^2} \sum_{b \neq a} \frac{e_b}{R_{ab}} B_{b|a} \hat{\mathbf{e}} \cdot [\mathbf{v}_a + \mathbf{n}_{ab}(\mathbf{v}_a \cdot \mathbf{n}_{ab})].$$

Since the shift in the interference pattern of charged quantum particles is given by the integral over time of the magnetic potential energy term of an appropriate Lagrangian, and since the Darwin Lagrangian is appropriate for this particular setup, this last term delivers the exact same shift. But in addition, this last term reveals that the effect is brought about by the magnetic fields generated by individual current-charge elements.

This account is gauge independent, and is local as well. It does not involve any gauge ontology and it brings the field center stage: after all, it has been the field that emanates from charge-current element what influences other charge-current elements all along.

Tim: Intrinsic Geometrical Structure and Fiber Bundles

The metaphysical moral of Tim Maudlin's talk was that relativistic quantum physics is hyper-local; that is, more local than classical physics. His view contrasts with that of Richard Healey, whose metaphysics is non-separable and therefore less local. The object of contention is of course the connection.

Healey argues that the connection cannot be used as representation of physical properties because there are distinct yet holonomy-equivalent connections (e.g. ω and ω') which may represent different distributions of properties for the same physical setup. The consequences of this ambivalence are empirical underdetermination and semantic inexpressibility, and combined with the epistemic inaccessibility of connections make the metaphysical picture that involves connections unfeasible.

The only reason for the connection's bad name, however, is over-counting, or so Maudlin argues. The real issue is whether there are many distinct connections that are consistent with uniquely determined physical parameters, like the magnetic flux in the Aharonov-Bohm effect, or just one. And the tension is resolved once we realize that the connection ought to be understood as a unique geometric object but with several (infinitely many, in fact) arithmetic representations or local representatives.

To clarify the confusion that has arisen from this conflation Maudlin used a series of examples. To that effect he first presented three one-dimensional spaces (R^1 , E^1 and V^1) and three two-dimensional spaces (R^2 , E^2 and V^2). A significant difference between numerical and geometric objects highlighted by these examples is that numerical objects (e.g. numbers that make up R^1 and R^2) are qualitative different from each other, whereas geometric objects (e.g. points that make up E^1 and E^2) are qualitatively identical. This difference implies that when we use numbers to represent geometric points there will be many ways of doing so and these ways will be mathematically distinct. The unassailable fact of the matter, however, is that the geometric objects are not as many as their arithmetic representations, and this is the first point where counting incorrectly may lead to confusion.

From spaces like these six several new spaces may be constructed but of interest here are product spaces and [more convoluted] fiber bundles. These spaces are of interest because of an attention-grabbing difference between them which has a bearing on metaphysical issues, namely the geometry of product spaces is completely determined by the geometry of their factors but the topological structure of [the other] fiber bundles is not.

When a product space $T \times S$ is employed to represent physical situations through, say, a function F , all the values of F lie in a common space, the space T , which attaches on each point of the physical or base space S . Since it is the same space T that attaches on each point of S , the values of F are drawn from the same space of possibilities and are therefore intrinsically comparable to each other. Each physical situation, then, picks a unique cross-section on the product space. This mathematical arrangement fosters a metaphysics of universals.

There maybe cases, however, where the spaces of physically possible states T_x that attach on the points x of the physical space S are generically like one another but there is no identity between them. In such cases it is not the notion of product space that

is appropriate but that of the [more convoluted] fiber bundle. T_x is called the fiber above x and there is no intrinsic relation among the points of different fibers (i.e. fibers above different points x), but the collection of all fibers is a topological space in its own right. In the neighborhood of any base space point x , the topology of the fiber bundle is isomorphic to that of the product space. Its overall topology though may be different, as the examples of the product bundle $E^1 \times S^1$ and the Möbius band demonstrate: the topology of the first is isomorphic to the topology of the cylinder, but the topology of the second is not due to the fact that it is non-orientable. Thus, starting from the same ingredients, the same fibers E^1 and the same base space S^1 , one may construct fiber bundles with different overall topologies.

This is the second point where the question of counting arises: how many fiber bundles of E^1 over S^1 ? The answer is “exactly 2”. The origin of the difference between the two fiber bundles appears to be a non-local feature associated with their global geometries: after all, locally the two fiber bundles look exactly the same. But sometimes appearances are misleading. There is no non-locality involved if we abide by the following criterion of locality: *divide the space into arbitrary overlapping neighborhoods, and specify the intrinsic structure of each neighborhood plus the overlaps. If the whole structure follows, it is “globally local”*. Thus, although it is globally and not locally that one is able to distinguish between the two fiber bundles, it is the totality of local facts that determines the overall structure uniquely and unequivocally and for this reason the difference is not non-local but hyper-local. This difference in the global structure of the two fiber bundles may be used in representing different physical situations.

Comparison of different points that lie in different fibers of a fiber bundle is effected by some additional structure, called the connection [;the connection ‘lives’ in the principal bundle whose fibers are the structure group of the theory]. The connection defines parallel transport of a vector in a fiber [that ‘lives’ in the associated fiber bundle] from one fiber to another. The vector is parallel-transported via a continuous path in the base space. If the connection is flat, the parallel transport of the vector from the fiber over a base-space-point p to the fiber over another base-space-point q via any continuous path will leave the vector unaffected. Now, if the [base] space is not simply connected but its curvature is locally zero everywhere, the representatives of the connection in its simply-connected neighborhoods that are localized around a point p can still be decomposed into flat connections; and there will be no difference on a vector that is parallel transported from point p to point q of a simply connected neighborhood regardless of the path. That is to say, so long as vectors are parallel transported in such neighborhoods alone no difference will be recorded. But once we put together all the patches that together go around the non-simply connected region of the base space (the one that is responsible for its multiple-connectedness), a difference will be recorded: the original vector and the transported one will point in different directions, so to speak.

The difference between the flat connection of a fiber bundle with simply connected base-space (which results in no change in parallel transported vectors) and the flat connection of a fiber bundle with non-simply-connected base space (which results in an overall change in parallel transported vectors) corresponds to the difference between having zero and non-zero flux inside an Aharonov-Bohm solenoid. This difference can be

used to resolve the apparent paradox associated with the effect: the connection when the flux is zero inside the solenoid is also zero, but the connection when there is flux is not zero inside is not zero outside either. The fact that the numerical value of the local representatives of the connections (i.e. the gauge potentials) may be chosen differently from point to point (due to gauge freedom) does not diminish the significance of there being only one connection per physical setup. The “many gauges” that can be chosen in each case are merely the distinct mathematical representations of a unique geometric object; an object whose uniqueness is affirmed by the fact that the shift of the interference pattern is independent of these choices. And this is the third point where counting correctly matters.

The fact that there may be more than one way of representing geometric objects with numbers is a source of ‘gauge freedom’ that is totally unproblematic. There is more than one arithmetic way of representing a particular geometric point, but no one asserts that the geometric object in question is not unique. Similarly, there is more than one way of representing arithmetically the connection of a particular physical setup at each base space point, but this does not mean that there is more than one connection.

The only thing that takes a blow from the correct counting of connections is the metaphysics of universals. A metaphysics of universals, which is compatible with product spaces, requires that for a given physical situation there is a uniquely determined correspondence between the points in base space and the magnitudes of a physical object in a space of universals (i.e. the space of the possible magnitudes of that object); even when this correspondence is expressed arithmetically in various ways. This amounts to choosing a cross-section in the product space which is unique for a particular physical situation, despite the ‘gauge freedom’ afforded by its multiple numerical representations. When the physical objects are represented by fiber bundles, though, there is another source of gauge freedom which undermines this kind of metaphysics. In this case there is not a single geometric space, or fiber, of possible magnitudes of the connection but multiple such spaces, or fibers, one for each base space point. But in any event the fiber bundle is unique, and so is the connection.

Antigone’s concern #1 (provided she understood Tim correctly):

The second kind of gauge freedom doesn’t have to do with the choice of the fiber-space. The fibers of a fiber bundle are fixed –e.g. for EM they are $U(1)$. Gauge freedom has to do with the fact that these fibers can be mapped onto themselves without changing anything physical. Gauge freedom has to do with the fact that we can pick not one but infinitely many cross sections in the principal fiber bundle of the potentials. Put differently, it has to do with the fact that we are not forced to pick a unique cross section for a physical situation because we can pick local representatives of the potentials *almost* as we please. For example: in the AB case we can choose a value of the potential locally; and another value further down the path; our only restriction in this case is that we cannot choose a cross section that is zero. This does not mean that we choose from different spaces of universals; it only means that we are relatively free to choose from within a single space, whereas in the product space case the value of a function for a particular physical situation is fixed. So, I’m not sure that the metaphysics of FBs is not compatible with a metaphysics of universals.

Antigone's concern #2:

Tim talked about holonomies briefly and said that connections contain more information than holonomies and that this is the reason they appear to be non-separable. My concern here is that Richard's cautious epistemologist will object that this extra information is empirically inaccessible and thus meaningless. I have to say that I am on the side of the prudent epistemologist.

He also said that holonomies will only tell you that the connections in the different neighborhoods are flat. This, I think, is not quite accurate. Knowledge of all the holonomies from a given base space point is sufficient for deriving all the gauge invariant content of the theory. In an AB setup the information you get from holonomies will tell you that the fields are non-zero somewhere and, therefore, that your connection cannot be zero everywhere.

Ward: Spontaneous Symmetry Breaking and the Higgs Mechanism

Following the lead of John Earman, Ward Struyve searched for the gauge invariant content of field theories that undergo spontaneous symmetry breaking (SSB) and the Higgs mechanism. The purpose of his exploration was to answer two questions posed by John Earman: “What does it mean to break a gauge symmetry?” and “If gauge is merely “descriptive fluff”, how can it have any physical consequences to break it?” Struyve addressed these questions by considering gauge invariant formulations of the classical theory, in both the Lagrangian and Hamiltonian picture.

SSB occurs when the Lagrangian density describing a system, e.g. a complex scalar field $\varphi = \varphi_1 + i\varphi_2$, is symmetric with respect to some symmetry group while the states of minimum energy (vacuum states) are not. An example of a Lagrangian density describing such a system is $\mathcal{L} = \frac{1}{2} \partial_\mu \varphi \partial^\mu \varphi - \left(\frac{m}{2} \varphi^2 + \frac{\lambda}{4} \varphi^4 \right) = \frac{1}{2} \partial_\mu \varphi \partial^\mu \varphi - V(\varphi)$, which is symmetric under global U(1) symmetry. When $m^2 > 0$, the vacuum state of the system (given by the minimum of V) is $\varphi = 0$ and unique, so that there is no SSB. When $m^2 < 0$ however, the energy has a local maximum at $\varphi = 0$ and the minima, which lie along the circle $|\varphi|^2 = -\frac{2m^2}{\lambda}$, form a set of degenerate vacua that are related to each other by a global U(1) transformation. Low energy fields are described by perturbations around these vacua. By choosing one particular vacuum and considering perturbations around that vacuum the original global U(1) symmetry is broken or hidden. The effective Lagrangian describing these perturbations (which is obtained by plugging the perturbation expansion in the original Lagrangian) concerns two new fields replacing the original fields: a massive scalar field and a massless scalar field (the so-called Goldstone boson).

The Higgs mechanism comes into play when the symmetry is a local U(1) symmetry. Contrary to the global U(1) symmetry, this symmetry implies indeterminism and hence gauge freedom. The Lagrangian is obtained by introducing a massless U(1) gauge field and replacing the partial derivative above by the covariant derivative:

$\mathcal{L} = \frac{1}{2} D_\mu \varphi D^\mu \varphi - \left(\frac{m}{2} \varphi^2 + \frac{\lambda}{4} \varphi^4 \right) - \frac{1}{4} F_{\mu\nu} F^{\mu\nu}$. The vacuum is still degenerate, but the minima are now connected by a local U(1) symmetry. If the aforementioned process of breaking or hiding the symmetry is repeated, the three original fields (the massive φ_1 , the massive φ_2 and the massless U(1) gauge field A_μ) are rearranged to give what looks like a massive φ_1 , a massless φ_2 and a massive photon, but also a peculiar mix of A_μ with φ_2 . At this point one usually “gauges away” the field φ_2 by employing the “unitary gauge”. The result is that A_μ becomes massive. One says that A_μ has gained mass by eating the field φ_2 .

Gauge freedom is associated with indeterminism and the restoration of determinism with gauge independent descriptions. There are at least two strategies one might adopt in order to get such a description. One strategy is to perform a field transformation that separates gauge independent variables or observables, whose motion is completely deterministic, from pure gauge degrees of freedom. One then ignores the latter degrees of freedom and works solely with the gauge independent variables. These variables could

potentially be regarded as physically real. A second strategy is gauge fixing, i.e. to impose extra conditions on the fields. It was argued that the gauged-fixed variables are in close connection to particular gauge independent variables that can be obtained by the first strategy. So potential ontologies arising from these strategies are but slightly different. This was illustrated for the unitary gauge. The corresponding application of the first strategy goes as follows.

Starting from the polar decomposition $\varphi = \varphi_1 + i\varphi_2 = \rho e^{i\theta} / \sqrt{2}$, with $\rho = \sqrt{2\varphi^* \varphi}$, $\theta = \ln(\varphi/\varphi^*)/2i$, a transformation is obtained, $(\varphi, A_\mu) \rightarrow (\tilde{\varphi}, \tilde{e}, A'_\mu)$, such that a local U(1) transformation $\varphi \rightarrow e^{i\alpha} \varphi$, $A_\mu \rightarrow A_\mu - \partial_\mu \alpha / e$ reduces to:

$$\begin{aligned} \rho &\rightarrow \rho \\ \theta &\rightarrow \theta + \alpha \\ A'_\mu &\rightarrow A'_\mu \end{aligned}$$

So ρ and A'_μ are observables, whereas θ is a pure gauge degree of freedom. The resulting Lagrangian for the observables ρ and A'_μ does no longer have a degenerate vacuum. By considering perturbations around the vacuum we get the same result as before: a massive scalar field and a massive photon. What's most important is that by following this strategy we get a clear picture of the observables involved without the intermediate cannibalistic step, where "descriptive fluff" must be eaten in order for physical degrees of freedom to be born.

Turning to the Hamiltonian treatment of the same system, Struyve presented two different ways of applying the first strategy. The first way uses similar variables as those just used in the Lagrangian picture. However, because these variables arise from the polar decomposition of φ one has to assume $\varphi \neq 0$. The second way concerns a transformation that separates out the observables that is defined over the complete field space. In this case there is still a residual symmetry for the gauge independent degrees of freedom, namely a global U(1) symmetry. This symmetry can be broken just as in the case of the free scalar field. The result is again the same as before: a massive scalar field and a massive photon. The reason a global U(1) symmetry did not show up in the first way of implementing the first strategy is that the field transformation could not be applied globally.

Antigone's critical comment on James' argument.

James's idea, namely that the AB effect is brought about by the magnetic field components which are produced by (multiple) charge-current elements, rests on his understanding of the content of the last term in the Darwin Lagrangian. He argues that the

term $\frac{e_a}{2c^2} \sum_{b \neq a} \frac{e_b}{R_{ab}} [\mathbf{v}_a \cdot \mathbf{v}_b + (\mathbf{v}_a \cdot \mathbf{n}_{ab})(\mathbf{v}_b \cdot \mathbf{n}_{ab})]$, which is the magnetic interaction energy between charge a and the charges b circulating in the solenoid, may be re-written so that the magnetic component field appears explicitly in it. This can be done if one takes into account the fact that $\frac{e_b}{R_{ab}} \mathbf{V}_b = B_{b|a} \frac{R_{ab} c \hat{\mathbf{v}}_b}{\sin(\theta_{\mathbf{R}_{ab}, \mathbf{v}_b})} \stackrel{def}{=} B_{b|a} \hat{\mathbf{e}}$, and substitutes $\frac{e_b}{R_{ab}} \mathbf{V}_b$ with $B_{b|a} \hat{\mathbf{e}}$.

Then the last term of the Lagrangian becomes $\frac{e_a}{2c^2} \sum_{b \neq a} \frac{e_b}{R_{ab}} B_{b|a} \hat{\mathbf{e}} \cdot [\mathbf{v}_a + \mathbf{n}_{ab}(\mathbf{v}_a \cdot \mathbf{n}_{ab})]$ and the magnetic component fields figure in it explicitly.

The problem I have with this exposition is the following. What appears explicitly in this term after the algebraic manipulation is *not* the magnetic component field but the magnitude of it. The magnetic component field that is produced by a charge b where a is given by the expression

$$\mathbf{B}_{b|a}(R_{ab}) = \left[\frac{q}{R_{ab}^2} \frac{v_b}{c} \mathbf{n}_{ab} \times \hat{\mathbf{v}}_b \right]_{ret} \quad (1)$$

and it is always perpendicular to the velocity \mathbf{V}_b of charge b . The expression that

replaces $\frac{e_b}{R_{ab}} \mathbf{V}_b$, however, is obviously parallel to \mathbf{V}_b since $\hat{\mathbf{e}} \stackrel{def}{=} \frac{R_{ab} c \hat{\mathbf{v}}_b}{\sin(\theta_{\mathbf{R}_{ab}, \mathbf{v}_b})}$.

The fact that the last term of the Lagrangian has got nothing to do with the action of the magnetic component field $\mathbf{B}_{b|a}$ on a may be seen also from the following. The magnetic component field that is produced by b acts on an electric charge a according to Lorentz's magnetic force law

$$\mathbf{F}_{b|a} = \mathbf{v}_a \times \mathbf{B}_{b|a} \quad (2).$$

Substitution of (1) into (2) yields

$$\mathbf{F}_{b|a} = \mathbf{v}_a \times \mathbf{B}_{b|a} = \left[\frac{q}{cR_{ab}^2} \right]_{ret} \mathbf{v}_a \times \mathbf{n}_{ab} \times \mathbf{v}_b = \left[\frac{q}{cR_{ab}^2} \right]_{ret} [(\mathbf{v}_a \cdot \mathbf{v}_b) \mathbf{n}_{ab} - (\mathbf{v}_a \cdot \mathbf{n}_{ab}) \mathbf{v}_b] \quad (3).$$

Expression (2), or (3) for that matter, is definitely not part of the last term of the Lagrangian despite the presence of the magnitude of the magnetic field in it, and this should come as no surprise. Lorentz's magnetic force law is an integral part of the equations of motion and it cannot be read off directly from the Lagrangian; for, electromagnetic Lagrangians contain information about the potential energies of interactions and not about the forces that may be involved. On the other hand, the magnitude of the magnetic field from each charge b is bound to show up in the

expression of its potential one way or another¹, but the magnetic field itself will not. This means, however, that we cannot interpret the last term of the Darwin Lagrangian as describing the interaction between the charge a and the magnetic field from b ; for, whether an electric or magnetic force acting on a particle is zero or not, only the physical set up and the equations of motion can tell not the Lagrangian.

Turning to the magnetic component fields and their action, a qualitative assessment of the AB physical situation reveals the following. At any given point of its classical trajectory, electron a experiences force due to the magnetic component fields produced by each of the b charges that circulate inside the solenoid. In the limit where the radius r of the solenoid is much smaller than the distance R_{ab} between an electron b inside the solenoid and an electron a outside, $R_{ab} \ll R_{ab} + 2r$. This means that at any given instant, electron a experiences equal but opposite forces from every two b electrons that lie diametrically opposite each other on the extension of R_{ab} . This, in turn, means that as much as one of the two component forces shifts electron a in one direction, that much the other component force shifts it in the opposite direction and the net magnetic force on a is zero. Therefore, at least in the limit where $R_{ab} \ll R_{ab} + 2r$, the shift in the interference pattern of the a electrons cannot be explained by the action of magnetic component fields as it cannot be explained by the action of the bulk field either.

One may argue that the case I presented is a limiting one and that in order to prove my point I have to show it to be true in general. I disagree. For one, if James's interpretation were fundamental, it should be valid on all approximations; and I have shown that in this particular approximation it is not valid. For another, the Darwin Lagrangian is itself an approximation that ignores the radiation degrees of freedom; therefore no universal conclusions can be drawn from it anyway.

¹ For example, the potential for the original AB setup (Aharonov and Bohm, 1959) is given by the expression $A_\theta = \frac{\phi}{2\pi r}$, where the magnetic flux $\phi \propto B$; and the potential for an A-B effect with two solenoids (Stovicek, 1993) is given by the expression $\mathbf{A} = \left(\frac{\hbar c}{e}\right) \text{grad}(\alpha\varphi_A + \beta\phi_B)$, where α and β are again the magnetic fluxes through the two solenoids and are proportional to the respective magnetic field magnitudes.

Brief responses to Tim:

I would like to add my voice to Gordon's and, in concert with him, say that one must face the music if one searches for an interpretation that employs a particular formalism of the theory. This involves taking the principal bundle seriously, but also analyzing what Tim's claim –that the connection is unique- may mean.

In Tim's view the fiber bundle (FB) formalism can solve the epistemological and semantic problems highlighted by Richard, if it can be shown that gauge equivalent potentials represent the same connection, which ought to be unique to each physical situation. And, in Tim's view again, we can tackle the problem without resorting to the entire formalism, for we can isolate the relevant parts and discard the rest. Thus, we should restrict our attention to FBs whose fibers *represent possible physical states associated with points in the base space*.

In danger of misrepresenting Tim, I take this to mean that the FBs he is referring to are the associated vector bundles, the FBs where matter fields typically 'live'. Assuming understanding, I believe that for the uniqueness-of-connection argument to fly one would have to show that starting from the associated bundle the connection is uniquely determined. But in order to define parallel transport in the associated bundle, one first defines vertical and horizontal sub-spaces of the tangent space to the associated bundle (in a way similar to that described by Gordon, the difference being that his FB was a principal one). And the vertical and the horizontal sub-spaces of the tangent space are defined through the connection ω of that 'lives' in the principal bundle. (See Isham, pp. 267-8.)

Is this reference to principal bundles unavoidable? Not if one does not appeal to FBs at all. In this case, however, the connection will be identified with the local representative $A_\mu^a(x)$ and thus we are back in the vicious circle we've been trying to avoid. Note that in this case the epistemological argument does have an impact.

If one appeals to FBs, on the other hand, the aforementioned reference indicates that one must appeal to the principal bundle. And this should be expected since, by definition, any FB is associated with a principal bundle, and all the information about the associated bundle is contained in the principal bundle.

One might argue, of course, that it is precisely the connection ω in the principal bundle which is uniquely determined. After all, ω , as a geometric entity, is unique. The problem with this idea is that these entities 'live' in an abstract geometric space and their manifestations in physical space are not uniquely determined. Tim could assert that these manifestations are only the numerical representations of the unique real geometric object, and I have to say that I find the idea very appealing. But the problem with it is that the geometric object is not present in the physical space(-time) between the solenoid and the electron.

Regarding holonomies and the information one gets from them, I would like to stress that knowledge of all holonomies from a point x in base space/manifold suffices for determining the gauge invariant content of the configuration –this is what Giles (1981) showed. When the physical space is treated as though it is non-simply-connected (and it's got to be for the AB effect to exist at all –note that the non-simple-connectedness does not require an infinitely long solenoid), physicists often employ the universal

covering space that is simply connected, and it is in this space that there is no “there”; but the information about the magnetic flux, or the associated potentials, is still encoded in the boundary conditions. In a non-simply connected space with zero magnetic flux through the hole, there is no AB effect (see Bernindo and Inomata (1981)), and if one works in the universal covering space this is reflected in the boundary conditions again.

Finally, I would like to comment briefly on Tim’s response to John Earman’s second question. John is asking two things:

- (a) How to connect the phase factors to the magnetic flux, and
- (b) How the FB formalism provides a basis for the desired physical interpretation, if at all.

Regarding the first question, Tim’s intuitive response is on the right track. Schulman (*Techniques and Applications of Path Integration*) identifies the (homotopy) phase factor with the holonomy/Wilson loop for the following reason. In general, in a multiply connected space the propagator can be written (roughly) as

$$K(x, t; y, 0) = \sum_{\beta} \exp(in\delta) K_{\beta}(x, t; y, 0),$$

where the sum over β is the sum over different

homotopy classes. When there is a magnetic field present and the space is multiply connected, δ is identified as $\delta = \frac{e}{\hbar c} \oint A_{\mu} dx^{\mu}$. Thus the phase factors are connected to the magnetic flux. The identification is justified, I believe, by the structural similarity between the fundamental homotopy group associated with the non-simply connected space-time and the structure group $U(1)$ of the principal bundle of the theory: the two are homomorphic.

But like John, I don’t see how this kind of identification may be of assistance in the quest for the desired interpretation. From a theoretical/geometric perspective we have more surplus structure than before because we have added yet another principal bundle: that whose structure group, \mathcal{H} , is the fundamental homotopy group of a plane with a point removed.

Response to James:

My first point was that Lagrangians do not contain information about forces explicitly. The fact that we typically attribute potential energies to force-fields (e.g. gravitational or electromagnetic) does not alter this fact. Now, in the case of the AB effect, since 1959 the idea has been that there is magnetic effect on electrons but no magnetic force acting on them. One who believes otherwise has to *show* that there is non-zero force where the electrons are. But one cannot infer this force directly from the fact that the Lagrangian contains terms for magnetic potential energy (or terms for magnetic interaction, as some might call them). Nor can one infer that it is the magnetic field components –let alone the magnetic force components— that give rise to the potential energy terms, because all that appears in the Lagrangian is the magnitude of the field components and not the field components themselves. Moving charges produce magnetic fields that may act on other moving charges, but saying that these fields or the forces

associated with them give rise to the magnetic potential energy is a metaphysical claim that needs to be substantiated. This brings us to my second point.

The second point was simply this: if the magnetic component forces were fundamentally responsible for the AB effect (and therefore for the presence of the magnetic potential energy terms in the Lagrangian), then these forces should account for the effect in all approximations, especially in the approximations in which the effect is accounted for by the potentials. The potentials account for the effect in the approximation where $R_{ab} + 2r \cong R_{ab}$ (in fact, this is a typical approximation), and they account for it whatever the Lagrangian –whether Darwin or not. This approximation does not require $R_{ab} \rightarrow \infty$, and therefore it does make the magnetic component fields zero. The result in this approximation relies on the fact that the component fields from every two diametrically opposite (and therefore moving in opposite directions) electrons in the solenoid will produce equal but opposite forces on any test-electron that will cancel out. The phase shift calculated through the Darwin Lagrangian is still non-zero (the velocities of the diametrically opposite electrons are equal and opposite, hence the subtraction will always be non-zero); but the actual force on the electron from the magnetic component fields is nil in this approximation and therefore the AB effect cannot be attributed to them.

Tim Maudlin

Responses to Antigone Nounou's Concerns on his Talk

Tim's response to concern #1:

There is a disconnect between the particular mathematical structure that Antigone mentions here and the structures I was discussing. Antigone is concerned with a principal fiber bundle- i.e. a fiber bundle in which the fibers are group operators. This adds another level of mathematical structure beyond anything I had in mind, and one would have to have a long discussion from the foundations to understand the physical significance of these groups. To take the particular example, Antigone writes: "The fibers of a fiber bundle are fixed -e.g. for EM they are $U(1)$." I'm not at all sure how to understand this. Suppose, for example, I have two physical objects in a two-dimensional space- two scalene triangles, to make the example concrete- that are in different regions of the space. For each of these objects, there is a group of transformations I could apply to it- a group of rotations- that has the structure of $U(1)$. Now there are two points. First, the physical meaning of the group is parasitic on the existence of different possible states of the objects, in this case different possible orientations of the object. The fibers I was considering were fibers that represent different possible physical states associated with points in the base space, not a group of operations on those states. So I don't know how to make direct physical sense of a fiber bundle whose fibers are $U(1)$: we have to have a discussion of the physical ontology that is being presupposed. This is made more complicated in quantum theory, since we use a (mathematical) wavefunction to represent some sort of physical reality, but everyone agrees that, e.g. the overall phase of a wavefunction has no physical significance. So we want to think of wavefunctions (at least) projectively to get a mathematical object that could be in closer to 1-to-1 correspondence with possible physical states. So as not to cause confusion, let me denominate the physical state, or condition, or whatever, that the (mathematical) wavefunction represents the "quantum state" of a system. We all recognize that the relationship between wavefunctions and quantum states is not 1-to-1: this is a classical example of gauge freedom. Now I suppose that when Antigone thinks of the fibers in her fiber bundle as $U(1)$, she is thinking of operators that operate on the wavefunction. But how this relates to any physical ontology is (to me) multiply obscure. In any case, it was not what I was talking about at all.

Now a second point. Let's take the case above where we understand the ontology: there are two scalene triangles at different locations in a 2-dimensional physical space, and for each there is a set of rotations that can be applied, which form the group $U(1)$. OK, in both locations there is *generically* the same group. But even in this case, there may not be any *particular* identification of the elements of these groups with one another. For example, for each triangle, there are two distinct 90° rotations possible. We can try to designate one of these a "clockwise" rotation and the other a "counterclockwise" rotation, but these designations are arbitrary, like calling one direction on a Euclidean straight line the "positive" direction and the other the "negative" direction. If we imagine that the two triangles live in *disconnected* spaces, or in a *non-orientable* connected space, then there is

no “fact” about which 90° rotation of the first triangle corresponds to which of the second: this is exactly the point about connections that I have been trying to make. So in such a case, even though each triangle (i.e. each point in space) has a $U(1)$ group of rotations associated with it, the fibers are not “fixed” (to use Antigone’s term) in the sense that there is a unique 1-1 mapping from the elements of the group associated with one point to the elements of $U(1)$ associated with the other. If the triangles happen to be in a connected, orientable space, then there will be a unique way to associate these groups with each other. But the 1-1 association problem becomes much more severe in a higher-dimensional space: if the space is three-dimensional, and the rotation group is $O(3)$, then there will be a unique association of particular rotations at one point with particular rotations at another only if the space is connected, orientable, and flat. This brings us back to the connection of the tangent bundle with which we are all familiar.

I have been using the example of rotations of physical objects that live in the space (rather than transformations among physical states that are represented by a fiber (NOT A PRINCIPAL FIBER!)) because these examples are familiar and easy to understand. They show that the phrase “the fibers are fixed” seems to miss the point I have been trying to make. In a $U(1)$ bundle, the fibers are all *generically* the same (they have the group structure $U(1)$), but it does not follow that there is any identification among elements of different fibers. Further, focusing on these principal fiber bundles, where the fibers are groups, makes the physical ontology extremely obscure to me, which is why I was discussing, e.g. vector bundles. Richard wants to give the vector bundles a somehow less important status in understanding the physical ontology, but I just don’t understand the picture. The group of operators that make up the fibers of a principal fiber bundle have to operate on something, and we only have a chance to understand what they represent when we understand the thing they operate on, and what *it* represents.

I would also recommend that we drop all reference to “potentials”. Antigone writes of “infinitely many cross sections in the principal fiber bundle of the potentials”. But the classical potentials- i.e. the vector and scalar potentials- all suffer from various gauge freedoms, as we all know. *They* are certainly not in 1-1 correspondence with possible physical conditions. And there are “gauge transformations” that take us from one set of potentials to another set that represents the same physical state. So the phrase “principal fiber bundle of the potentials” suggests (at least to me) a fiber bundle whose fibers are groups of gauge transformations in the trivial sense: transformations from one set of potentials that represents a connection to another set of potentials *that represent the same connection*. Without further argumentation, these sorts of gauge transformations are of no concern whatever once we take the connection as the physical reality. Richard’s position, as I understand it, is that even once we move to talk of connections, there is still an underdetermination problem of some kind. The main thrust of my argument is that the connections are not underdetermined: e.g. given the flux through the solenoid, there is only one connection on the relevant bundle that is possible, which will be flat in the annulus outside the solenoid, but nonetheless different in that region when the flux is different. Maybe I’m wrong about this, but focusing on the potentials is just going to obscure the argument.

Tim's response to concern #2:

No- the connections contain more *local* information than the holonomies. If the argument goes through, the *totality* of the holonomies contains exactly the same information as the connection- i.e. there is only one connection compatible with the totality of the holonomies. I don't think the epistemological argument cuts any ice. The gauge transformations, on this view, are just that: transformations between distinct mathematical representations of the same physical state.

What I had in mind about the holonomies for the second point is this: consider only the fiber bundle over the annulus outside the solenoid. Now Richard made a claim about the supervenience of holonomies for big loops on holonomies for small loops: the claim was something like this: since you can make big loops out of small loops, if you fix the holonomies on all the small loops, you automatically fix the holonomies on the big ones. (At least, this is what I understood him to say.) In the case where restrict ourselves to the annulus, this is untrue. The small loops (none of which are big enough to circle the hole in the center of the annulus) will tell us no more and no less than that the connection is flat. They will contain no information about the holonomies of big loops that encircle the center hole. This can only happen because the base space is not simply connected.

Now even in the A-B effect, I don't think it is accurate to say that "the information you get from holonomies will tell you that the fields are non-zero somewhere", although this is subtle. In the actual A-B effect in the lab, the physical space is, of course, simply connected, so the totality of holonomies for the annulus do imply the existence of a non-zero field (i.e. a non-flat connection) somewhere in the hole in the middle of the annulus. But the effect itself does not require that the physical space be simply connected. Suppose physical space itself were not simply connected, so there is nothing "in the hole in the center of the annulus"- i.e. there is no "there" there. The physics could be just the same (I take it Richard agrees here!): particles whose wavefunctions go both ways around the annulus could display a phase shift, even though the connection is flat everywhere (no non-zero field anywhere). For me, it is a matter of the total connection, and different connections that are everywhere flat are possible (that's the point about the frustums of the cones). For Richard, as we have seen, the holonomies of loop that enclose the "hole" will not supervene on the holonomies of loops that do not. In either case, the effect could be produced in a non-simply-connected space-time without any non-zero fields anywhere.

General Comment & Reply to Richard

I would like to make one general comment about the A-B effect and our discussion of it. What we ultimately want, of course, is a complete physical account of the A-B effect. Such an account would have two components:

- 1) An account of why there are interference bands in the first place and
- 2) An account of why the bands shift when the current in the solenoid is changed.

Problem 1), of course, requires an interpretation of quantum theory, the status of the wavefunction, what are the local beables, etc. This problem can equally well be discussed for the plain vanilla two-slit experiment, and should be. At the conference, nobody took

up this problem, and Richard's offhand comment to the effect that he was no longer pursuing a modal interpretation but rather looking back to Bohr suggests that any such discussion would have been very complicated and contentious. I said myself that I was not going to take up the general interpretational problem. But we should clearly demarcate questions that cannot even begin to be addressed without attacking this problem from those that can. I imagine that general questions about a "causal account" of the phenomena cannot be addressed without attacking it.

Problem 2), however, can be discussed without completely settling problem 1)- or at least some remarks can be made about it. This is what I was undertaking to do. I take the general outline of the problem to be this:

- a) The observable phenomena (in this case the locations of the interference bands, but nevermind exactly what the phenomena are) change when the current in the solenoid changes.
- b) The things that produce the phenomena (the electrons- but nevermind exactly what *they* are) are shielded from entering the interior region where the solenoid is, so whatever the physical difference it is that accounts for the change, it must be *outside* the solenoid. (This rules out a certain kind of unmediated action-at-a-distance.)

The combination of a) and b) poses the question: what exactly *outside* the solenoid changes when the current changes? And what *doesn't* change?

Here are four answers:

Classical electromagnetics: The vector and scalar potentials change, while the E and B fields do no change.

Tim: The connection on a fiber bundle changes, while the curvature of that connection does not change.

Richard: The holonomies for loops that enclose the solenoid changes while the holonomies for loops that do not enclose it does not change.

James: The component magnetic fields at each point changes while the bulk field at each point does not change.

Now in order to endorse any of these answers as having the wherewithal to provide a physical explanation, one must ascribe some sort of *physical reality* to the thing that changes (if we rule out action-at-a-distance). So the problem with the Classical electromagnetics answer is that it requires ascribing physical reality to the potentials, and a physical difference to potentials that imply the same fields. This contradicts the classical understanding of the potentials, and raises questions about whether one must regard even gauge-equivalent potentials as different. No one wants to do this.

My solution solves this last problem if it is correct (as I argue) that gauge-equivalent potentials represent the same connection. I agree with Richard's claim: the holonomies for loops that do not enclose the solenoid are unchanged, and the holonomies for loops that do are changed. But I want to see the holonomies as derivative entities. Mathematically, one calculates the holonomies by evaluating intergrals around closed loops- and this is how I see them. But for those integrals to have any meaning, you need the connection. After all, integrating is a

prototypical *local* procedure. I can't see how to make sense of any of this on Richard's ontology, where the holonomies are somehow primitive and atomic.

Neither my solution nor Richard's has any need to make a distinction between bulk and component fields, so they both are in a different ballpark than James.

Now a response to Richard. First, as a small point, there is really no need to replace my cone with a torus- the set of directions that a vector on the cone can point has the structure of $U(1)$. So picturing this as a cone is OK.

Now Richard wants to argue that I'm wrong about the solenoid/connection relation: he wants to say that a certain state of the solenoid is compatible with only one set of holonomies but with multiple distinct connections. But although Richard has denied that his argument is just the hole argument redux, I think it is exactly the hole argument redux. Anyway, his argument about multiple "lifts" with the same holonomies is, I think, the hole argument redux, with his "change in connection" playing the role of the diffeomorphisms in the hole argument. Now I don't want to rehearse the interpretations of the hole argument again, but I really see no structural difference here at all. And I think that Richard is committed to saying what I would deny: that there are multiple flat connections for a fiber bundle over a simply-connected base space. So I recommend thinking about just this example. Can we agree, at least, that if there is only *one* such connection, that Richard's underdetermination argument fails? Because if his argument for multiplicity works at all, it works in this case.

A couple quick comments:

First, I think Richard and I are making progress- indeed, where he says he thinks I see an analogy to the hole problem is exactly where I do see an analogy to the hole problem! This also comes through when Gordon explicitly analogizes the symmetry group of the principle bundle to the set of diffeomorphisms in GR- the source of the hole problem. I have already made my stand on the hole problem, so there is no need to go over that ground again.

Second, if I have understood John's second question- "What I don't see is how to connect the phase factors to the magnetic flux."- I think is clear at least in outline in terms of the fiber bundle- assuming my analogy is really strong enough to go through. Consider, yet again, my humble cone. Considering only the frustum, one can say it is everywhere flat, and that from that fact alone one cannot determine the cone angle, and hence the angular defect (and hence the phase shift in the wavefunction) for a loop around the center. But if you consider the *whole* cone, smoothing out the cone point into some sort of rounded cap, then one *can* calculate the angular defect for a loop around the center. And the curvature of the cap *just is* the magnetic flux in the solenoid.

Lastly- and this is a very general comment- I do not see how to answer the more detailed ontological questions without a careful discussion of exactly what the physical degrees of freedom in the wavefunction are. I take it we all agree that mathematically different wavefunctions can represent the same physical situation- e.g. wavefunctions shifted by an overall phase. Until we sort this out, we can't really distinguish a mere gauge symmetry (connecting two mathematically different objects that represent the same physical state) from a physical symmetry (two different physical situations that are empirically equivalent, according to some standard of what can be observed). I don't think that fancier mathematics will help us with that question.

John Earman

(1) One of the aspects of gauge I would like to see explored is the connection between gauge and other concepts, e.g. superselection rules (SSRs). Define a gauge transformation in QM (or QFT) as a non-trivial unitary that commutes with the algebra of observables M (assumed here to be a von Neumann algebra). Then it follows immediately that a SSR is implicated because M acts reducibly on the Hilbert space, which breaks down into a direct sum of superselection sectors. There are some easy but profound implications. For example, if the gauge group is non-commutative then there is no “complete set of commuting observables” i.e. no maximal abelian subalgebra of M . This initially caused some consternation and led to some dubious arguments, e.g. the modus tollens move that says that since the existence of a complete set of commuting observables is essential to QM, it follows that in the case of a system of identical particles there can only be bosons or fermions.

(2) Naïve question on the topological explanation of the AB effect. Suppose that the wavefunction cannot penetrate the solenoid. Then the configuration space is non-simply connected. Quantize using this space. The resulting wave function can be multiple-valued. Argue that the possible phase factors must belong to one-dim unitary representations of the first homotopy group of the configuration space (this can be rigorously justified). There we have it—the mathematical form of the AB effect without any of the usual BS. What I don’t see is how to connect the phase factors to the magnetic flux. And I don’t see how fibre bundle formalism helps—it provides a means of representing the above but not (?) a basis for the desired physical interpretation.

(3) When to recognize gauge can have an important effect on the outcome. Simple example: (a) Start with the standard configuration space for a system of N identical classical particles. Then quantize using this configuration space. Afterwards take account of the gauge freedom by declaring that the obvious unitary representation of the permutation group $S(N)$ is a gauge symmetry in the sense of (1). (b) Alternatively recognize gauge freedom from the beginning using the configuration space obtained by quotienting the configuration space of (a) by $S(N)$ and then quantizing. In approach (a) we end up with bosons, fermions, and paraparticles. In approach (b) we end up with only bosons and fermions for space dim at least 3, and anyons for space dim 2.

(4) The point made in (3) is connected to some morals Tim Maudlin wanted us to draw. Start by postulating some entities (where entity is construed broadly, e.g. it might be spacetime) and then using these entities fashion some (putative) laws. The solution set of the laws inherits the symmetries of the postulated entities. For various reasons one might wish to treat a symmetry of the solution set as a gauge symmetry in the sense of relating equivalent descriptions of the same physical situation. If so, says Tim (if I understood him correctly), it is a legitimate challenge to ask how the initial postulation can be modified so as to remove the (alleged) descriptive redundancy that results in the gauge symmetry. And (if I understood Tim) if you can’t meet the challenge, then you should not be so quick to see gauge freedom. Without entirely endorsing Tim’s point of view, I

agree that it is important to raise the challenge since it leads to a realization that physically inequivalent quantizations can result.

Richard Healey

1. While I applaud Tim's emphasis on the importance of not being confused into mis-counting structures by the use of overly-sophisticated mathematics to represent them, I think he is still mis-counting fiber-bundle connections. At the workshop I tried to press him on what properties/tropes he thought were represented by the connection on a *principal* (not vector) fiber bundle, but I now think my objection may be made without an answer to this question. So concentrate on the *matter* whose interference is observed in the AB effect. Tim takes this to be represented by a vector bundle, so the fiber above each point is an element of a vector space. But since what determines the shift in the interference pattern is just the *phase* relations of the matter (field or wave-function), perhaps we can agree to focus instead on a fiber bundle in which the typical fiber is not a vector space but a space with the structure of a circle—corresponding to the complex numbers of modulus 1, representing the matter phase at a point.

For a geometric model we now have not a cone, but a torus: and the issue is how to think of a connection on this torus that tells one what points of circles “above” neighboring base points are themselves neighbors. The points at the bottom of the torus correspond to points of the base space (for example, a circle in physical space surrounding the solenoid in the static magnetic AB effect). A connection on the bundle tells you what counts as “same phase as” as one moves around in the torus, by defining a notion of horizontal lift of a smooth curve in the base space to a curve on the torus: in this simple case the only relevant closed curves are those that trace out the base space circle some integral number of times. So how many connections are there?

An infinite number. Start from any point in the torus above base point m , and trace out an arbitrary smooth curve returning to a point in the fiber above p : that defines a connection. But these partition into equivalence classes, as follows. Consider any smooth closed curve C once around the torus from point e above m : this returns to point f above m making an angle 2 with e around the fiber above m . This defines a connection A : points above neighboring points of the base space are themselves neighbors just in case the tangent to a curve $C(N)$ “parallel” to C (i.e. displaced from C by the same fiber angle N above every point) points from one to the other. Now consider a *different* smooth curve C^* that goes once around the torus from point e above m and *also* returns to point f above m . This defines a connection B in a similar manner. But while A, B are distinct connections, they have the same holonomies, since going once around the torus moves you through the same fiber angle 2 , no matter whether the connection is A or B . A, B lie in the same equivalence class. There is an infinite number of such *classes*, each corresponding to a different angle 2 : each class itself contains an infinite number of holonomy-equivalent connections.

The angle of a cone in Tim's model corresponds to the angle 2 in my torus model. It picks out not a single connection, but a holonomy-equivalence class of connections. The way to see this geometrically is to think of the cone as composed of (infinitesimally) thin straight lines, and sliding each of these parallel to itself up and down on the cone, so as to turn one closed smooth curve around the cone into another. That will change the connection but not the holonomy, which remains fixed by the angle of the cone.

(Here it helps to remember that each line has the structure of a *complex* vector space, and the only permissible “slidings” change the phase, not the modulus, of the vector. That’s one reason why I think of the phase bundle as a better geometric model. In it, a change of connection within a holonomy-equivalence class corresponds to rotating the wheels of an unnumbered toroidal “combination lock” with continuously many rings so as to preserve the overall difference in the angle as you go once around the lock.)

What does all this amount to physically? The only phase relations that are physically comparable (experimentally or theoretically) are phase relations *at a point*: that is where interference phenomena manifest themselves, and phase relations are manifested only in interference phenomena. Interference phenomena can distinguish between *equivalence classes* of connections, since each equivalence class has a corresponding holonomy. But they cannot discriminate among holonomy-equivalent connections.

Now make the torus model a bit more realistic by making the base space into an annulus instead of a circle, and considering the phase bundle above this—geometrically, a circle “above” every annulus point, with a bundle connection formed by smoothly joining the connections on all the corresponding tori. Cover the annulus by a finite number (say, three) “patches” in the base space none of which includes the region near the center of the circle. What kinds of comparisons are possible within each patch, and between patches? Comparisons of holonomies around closed curves within each patch are possible: they may be operationally defined by idealized interference experiments within each patch. Such comparisons will not permit one to tell what constant current, if any, is flowing through the central solenoid: all the holonomies will have the same value—1. Moreover, these operations don’t permit one to single out a connection within a patch from the holonomy-equivalence class ideally revealed by those interference experiments. Parallel transport within a patch has no operational correlate, since phases can be compared only at a point. Moreover, the theoretical description of goings on within a patch is itself insensitive to a choice of connection from within the holonomy-equivalence class. That follows from gauge invariance.

What comparisons can be made in regions where patches overlap? Again, only holonomies around closed curves in such overlap regions are accessible to both overlapping regions. So no strictly *local* comparison or matching is possible between overlapping regions (though they can compare, and will agree on, the holonomies for arbitrarily small closed curves in the overlap region). Moreover, exchanging information about the results of internal comparisons between regions still leaves out information about holonomies of curves that pass through both regions, though theory warrants an inference to what those holonomies must be, given the “internal” holonomies, provided the union of the regions is simply-connected. When we consider *all* the regions covering the annulus, the results of pairwise comparisons do not suffice to determine the holonomies of curves circling the center, even when occupants of the regions are permitted to share the results of their individual comparisons. A region has no access to the connection within that region: *a fortiori* no comparison of connections is possible between regions. And there is no operational or theoretical content to an attempt to distinguish among holonomy-equivalent connections in the entire annulus.

2. Several questions and discussions at the workshop prompted me to think more about how the Holonomy Interpretation of classical EM handles locality. One thing to emerge

is that Local Action must be applicable even in situations in which loops L, M are neither spacelike nor null separated. Understood as requiring all distant influences to be *mediated*, it must at least require such mediation between *timelike* separated loops through a continuous sequence of intermediate loops. But that is still not enough.

Consider, for example, two similar circular loops that each occupy the same region of space in some frame, but lie on different simultaneity hyperplanes of that frame. If these hyperplanes are separated by a short enough time interval in that frame, the loops will be neither spacelike, null, nor timelike separated. But Local Action should still require mediation from the earlier to the later of an influence on properties of the earlier.

Maybe we can appeal to the general notions of the *causal past* and *causal future* of a loop, where the causal past (future) of a loop is the union of the past(future) light cones of all the points on the loop. Relativistic Locality would require that a loop lie in the causal future of all causes of its holonomy properties: while Local Action would require that for an influence on holonomy properties of loop L to have an effect on those of loop M in its causal future F_L , F_L must contain a continuous intermediate sequence of loops with holonomy properties nomologically related to those of L and M (provided there are no nomologically independent interfering influences in F_L).

Note that the loops referred to in these reformulations need not be spacelike, and nor need a pair of loops such as the L, M mentioned in Local Action be timelike separated (though they clearly cannot be spacelike separated without violating Relativistic Locality).

The claim now would be that the Holonomy Interpretation satisfies both Relativistic Locality and

Local Action, as so understood. Why should one accept this claim?

At the workshop, I admitted to feeling chagrined by the need to appeal to a particular gauge in arguing that holonomy properties propagate in accordance with Relativistic Locality. I shouldn't have! While the argument goes through most simply using a preferred gauge (that corresponding to what James Mattingly once called the current-field), it cannot depend on a choice of gauge, since the premises and conclusion don't, and parallel reasoning would be valid (but more messy) with any other choice of gauge. In this context, choosing a gauge is analogous to choosing a particular model out of a diffeomorphically-equivalent set in a generally covariant theory. One shouldn't feed chagrined when appealing to a model with a nicely symmetric metric in general relativity to prove a diffeomorphism-invariant result.

3. Some further thoughts on what's wrong with the Localized Properties interpretation of a classical YM gauge field. Here is how I stated the Localized property Interpretation (LI):

ω directly represents gauge potential properties: so holonomy-equivalent connections ω , ω^* may be used to represent different distributions of gauge potential properties over M

Perhaps this does not quite correctly state the view I wish to argue against, which includes a stronger claim, namely:

ω , ω^* may be used to represent different distributions of gauge potential properties over M in such a way as to leave open the question as to whether ω , ω^* (or indeed some other

holonomy-equivalent connection) correctly represents the actual distribution of gauge potential properties over M

Assuming there *are* localized gauge potential properties distributed over space-time points in a way compatible with the actual holonomy properties on space-time loops, one can simply decide to employ a system of representation under which these are represented by ω rather than ω^* .

Oliver Pooley convinced me that, having made that choice, it is legitimate to apply that *same* system of representation to ω^* to represent a *different* (but nonactual) distribution of localized gauge potential properties. But of course having done so it is not an empirical question whether the actual distribution is correctly represented by ω rather than ω^* : ω wins by default.

So even though ω , ω^* may (in this way) be used to represent different distributions of gauge potential properties over M , they may *not* be used to represent different distributions of gauge potential properties over M in such a way as to leave open the question as to whether ω , ω^* (or indeed some other holonomy-equivalent connection) correctly represents the actual distribution of gauge potential properties over M .

That is why one cannot even raise an empirical question as to which (if either) of ω , ω^* correctly represents the actual distribution of gauge potential properties over M . And if one cannot even raise this empirical question, then of course no experiment or observation can answer it.

Reply to Tim's response

Good: we are making progress!

Here are a couple of reasons to prefer the torus to the cone. A smooth relative rotation of neighboring infinitesimal straight pieces on the cone's surface, all around the cone, may involve radical changes in its geometry, so that it is no longer anything like the original cone: it may have to twist wildly into a ribbon shape in a way that would require *stretching* as well as maybe one or more rotations of some of these pieces "end over end". The surface would not stay flat: just try it with a paper cone! Also, a relative rotation of neighboring straight pieces would leave a single fixed point where they touch both before and after the rotation—where the axis of the rotation intersects the surface. A smooth choice of such fixed points around the cone effects a choice of connection from among the holonomy-equivalent connections that each give the same angular deficit after a single circuit of the cone's surface.

Since I don't want to allow transformations that change the curvature, I'll stick with my torus, representing not the magnitude but (at most) the phase of a complex vector representing the wave-function/field at the base point "below".

I do claim that there are multiple connections on this phase bundle over a simply connected base space, even when all horizontal lifts of closed curves in the base space themselves close in the bundle, in which case the bundle curvature is zero (as is the curvature of the complex vector bundle representing the wave-function's magnitude as well as phase).

N.B. Don't confuse the curvature of the torus's surface with the curvature of the phase bundle it depicts! The torus surface's curvature derives from its metric: the phase bundle's does not.

Each of an infinite class of distinct bundle connections is compatible with all holonomies of closed curves in the base space having the same value 1. Curvature is characterizable in terms of holonomies/holonomy-equivalence classes of connections. There are an infinite number of distinct connections compatible with any given curvature, whether or not the curvature is zero.

What does any of these connections represent, physically? The answer to that question depends on both us and the world. **If** the phase at a point represented some physical feature localized at that point, **then** by *fiat* we could choose to use connection ω to represent the actual distribution of such features over points (at this stage this just amounts to a *labeling* system for the points): but we could equally well have chosen ω^* instead. Having chosen to represent the actual distribution by ω , we could then use a *distinct* connection ω^* to represent a *different*, and so counterfactual, distribution by adopting the same mode of representation we adopted for ω . But then we would know *a priori* that the actual distribution is not represented by ω^* but by ω (maybe this is why Tim is reminded of the hole argument, where his metric essentialism depends on a similar move). What we *can't* do is give rich enough content to a mode of representation shared by ω , ω^* to be able to raise an empirical question, given that mode of representation, as to which (if either) *correctly* represents the actual distribution of localized phase properties. Since we can't raise such an empirical question, we can't answer it by experiments.

That is why I think we should conclude that *there are no* localized phase properties to be

represented---there are only non-separable holonomy properties. We can represent these directly by holonomies, less directly by an arbitrarily chose connection from the empirically appropriate holonomy-equivalence class, or even less directly by a gauge potential (or equivalence class of gauge potentials) on the base manifold (absent topological complications).

James Mattingly

Antigone seems to be making two distinct objections here. One arises out of an interpretive stance---that I am not understanding correctly the role that the magnetic field strength expression plays in the Darwin Lagrangian---and the other is a claim about the empirical facts---that the component fields I write down cannot have the effect I claim for them on any understanding of the role that the magnetic field strength expression plays.

Let me begin with the latter. I think this is easily cleared up because the objection is based on incorrect analysis of the physics. Antigone claims that the force on the exterior electrons from magnetic fields arising from any pair of moving charges on opposite sides of the solenoid should cancel in the limit of small solenoid and large distances from the solenoid to the exterior electrons. The intuitive method she uses to take this limit has thrown her off. The analysis of the magnetic field far from a current loop goes like the square of the radius of the loop. It is true that if we really take the limit of that radius vanishing then there is no magnetic field far from the loop. On the other hand, there is also no magnetic field at the loop itself. And indeed there is no loop. In order to maintain a finite current in this limit we need the electrons that comprise it to be moving infinitely fast. I think Antigone is quite right to point out that the net magnetic field due to a current loop arises from the slight differences in the distances from members of electron pairs to the observation point. It is not, however, correct to assume that these differences become negligible in the far-field limit. The best thing to do, I suggest, is just to calculate the terms as they are given in my proposal.

I'd like to be clear about this point. The calculation I perform that sits at the base of my proposal is pretty easy, and it shows clearly that the proper magnitude of the phase is recovered on integrating around any loop around the solenoid. This fact is made clear by the series of equations

$$\begin{aligned}
 \Delta\phi &= \Delta\phi_{path1} - \Delta\phi_{path2} = \int_0^t (L_{path1} - L_{path2}) dt \\
 &= \int_0^t \frac{e}{2c^2} \left[\sum_{b \neq e_1} B_b \boldsymbol{\kappa} \cdot [\mathbf{v}_{e_1} + \mathbf{n}_{be_1} (\mathbf{v}_{e_1} \cdot \mathbf{n}_{be_1})] - \sum_{b \neq e_2} B_b \boldsymbol{\kappa} \cdot [\mathbf{v}_{e_2} + \mathbf{n}_{be_2} (\mathbf{v}_{e_2} \cdot \mathbf{n}_{be_2})] \right] dt' \\
 &= \frac{e}{c} \int_0^t (\mathbf{A}_1 \cdot \mathbf{v}_1 - \mathbf{A}_2 \cdot \mathbf{v}_2) dt' = \frac{e}{c} \int_0^x (\mathbf{A}_1 \cdot d\mathbf{x}_1 - \mathbf{A}_2 \cdot d\mathbf{x}_2) = \frac{e}{c} \oint \mathbf{A} \cdot d\mathbf{x}
 \end{aligned} \tag{1}$$

As far as the mathematics is concerned, this proposal is identical to the standard account in terms of the change in phase that arises from the (ill-defined) vector potential.

So now about the interpretive point. Antigone thinks that because the magnetic force term does not arise in the Lagrangian then it cannot be the action of the magnetic field that is causing the change of phase. I don't get it. The force term for Newtonian gravity

does not appear explicitly in the gravitational Lagrangian, but we do typically attribute the gravitational potential energy to the action of the gravitational field. So it can't be as easy as saying that the interpretation I offer is illegitimate just because we cannot see explicitly the force term in the Lagrangian. I just can't understand the point Antigone is making about the origins of these expressions.

I'm happy to grant that force terms do not appear explicitly in Lagrangians. But it is not at all obvious to me that we should not think about the potential energy terms in Lagrangians by reference to the interaction forces that give rise to these potentials. In fact, it is precisely that assumption that I am trying to challenge with an *explicit* prescription for how to do so.

James on Ward and again on Antigone

- Ward raises an important point, that there is still gauge freedom in the wave function itself that is not dispelled by moving away from a formulation of the theory based on potentials. I don't have a complete answer here, but here is a suggestion.

First the question of gauge freedom for electrodynamics arises in classical, semi-classical, and quantum electrodynamics. What I am looking for is an interpretation of the electromagnetic fields that displays what is common across all three of these. Of course there might not be much, and so the project may well fail. But as my point of departure I focus first on the case of the electromagnetic fields and ask whether it is necessary to introduce essentially gauge-dependent quantities in accounting for their efficacy. I think the answer is clearly no in the case of the classical theory. And I think I've given reasons for thinking the answer is no in the semiclassical theory. These reasons detach from questions of interpretation of the quantum wave function in the sense that they piggy-back naturally atop standard accounts that concern themselves only with the change in phase around a loop that arises from some given choice of gauge. In part what I have done is show how to take one of these choices of gauge and rewrite it terms of local facts about the electromagnetic field strengths generated by the current elements. What my approach shows is that one can avoid picking a gauge, because the account is not in terms of potentials, while also appealing only to properties of spacetime points rather than regions or loops. And so the strictly electro-dynamical gauge freedom seems to have been eliminated.

And yet the fact remains that there is remaining gauge freedom in the wave function itself. But now we can at least suppose that that freedom arises not from features of the electromagnetic fields, but rather from facts about how we must represent wave functions in our theory. We know that local gauge transformations of the wave function must be accompanied by subtracting the derivative of that transformation from the covariant derivative. But on my account we have no ground for assuming that the connection between the covariant derivative and these gauge transformations. So I think we are free to regard the gauge freedom that remains in the theory as a representational issue, and one that has little to tell us about the ontology of the theory.

- I'm sorry to say that Antigone and I are not making much progress. I don't know what to say about her insistence that the magnetic field of a loop of current vanishes in the limit of small loops and large distances to the observation point. The calculation showing that the field goes like the square of the radius of the loop and the inverse cube of the observation distance is pretty straightforward, and appears, e.g., in section 5.5 of Jackson.

The mistake Antigone is making here, perhaps, is to make the $R_{ab} + 2r \cong R_{ab}$ approximation a little too early in her calculation. But in any case, a moving charge far from a loop of current does indeed experience a Lorentz force.

And in any case, I think Antigone is missing the point of the argument I am making. I do not claim that any overall force would be acting on a classical particle traveling about the solenoid. Indeed I think there would not be such a force. (But only because of the presence of all the other loops of current composing the solenoid.) Instead I claim that one can use the magnetic field that $\{em\}$ would arise from the presence of a single element current at some observation point to calculate what might be thought of as the phase change at that point induced by that element of current. And then one can add up the corresponding influences of all the other elements of current at that observation point to calculate what might be thought of as the total phase change at that point due to the presence of the charge-current elements in the causal past of the observation point. And then, whatever one wants to say about the character of the wave function in quantum mechanics, and however one wants to interpret what it is that is being acted on when electrons are scattered around the solenoid in a diffraction experiment, one can say at least the following: Calculating the fringe shift in the diffraction pattern using as argument to the action integral the localized property defined above---i.e., the phase change due to charge-currents in the causal past---yields the correct result. This result is obtained without making reference to anything like a vector potential, a holonomy, a hoop group, the frustum of a cone or anything else like that. It is obtained by instead taking note of the way that charge-currents give rise to electromagnetic fields.

All of that is trivially correct, and should be obvious. The real issue is how one ought to interpret the fact that such a construction is available. I tried to make a case at the workshop for an interpretation that takes seriously the ontological robustness of the component fields. Perhaps that case can be made stronger with a more compelling story about the ontology of the wave function itself. I plan to give it a shot.

Ward Struyve

Concerning James' approach to the A-B effect

James claims to have a gauge independent and local account of the A-B effect. However, I don't think the account is completely gauge independent. The A-B effect concerns quantum particles, described by a wavefunction, coupled to an electromagnetic field, described by a vector potential. Usually it is understood that a gauge transformation involves a phase transformation of the wavefunction, together with an associated transformation of the vector potential. In James' account a gauge independent ontology is introduced for the electromagnetic field (namely the component field), but apparently not for the quantum particles (if I understood correctly, the phase of the wavefunction is taken to be ontologically real). So with the usual understanding of gauge symmetry, the proposed ontology does not seem completely gauge independent. Is it possible to introduce a gauge independent ontology for the quantum particles too such that the account of the A-B remains local?

Comments by Gordon Belot

Philosophers of physics sometimes seem to think that principal fibre bundles are particularly horrible things. This is going too far. Principal fibre bundles are *somewhat* horrible things—but they are in fact just manifolds carrying nice group actions, and you can't read much philosophy of physics without coming across all of the notions needed to characterize principal fibre bundles.

Here are the details.

- A *Lie group* G is a group whose elements have been equipped with a manifold structure, in such a way that the operations of multiplication and the taking of inverses are given by smooth maps.
- Let G be a group and X be a set. Let φ be a function from $G \times X$ to X ; we usually write $g \cdot x$ for $\varphi(g, x)$. φ is an *action* of G on X if: (i) $g \cdot (h \cdot x) = (gh) \cdot x$ (for all g and h in G and x in X); and (ii) $e \cdot x = x$ (for all x in X with e the identity in G).
- An action is *free* if $g \cdot x = x$ implies $g = e$ (no element of G other than the identity fixes any point in X).
- Let G be a compact Lie group. A *principal G -bundle* is a manifold P equipped with a smooth, free G -action.

And that is all. It may be helpful to add two sets of remarks: one set providing a commentary on the above definitions; the other sketching the central role of principal fibre bundles in Yang–Mills theories.

Commentary on these definitions.

- In the first instance, one is typically interested in Lie groups whose elements are real or complex matrices. Let's focus on the real case. Let $E(n^2)$ be the space of n^2 -tuples of real numbers equipped with the obvious Euclidean metric. Each $n \times n$ real matrix corresponds to a point in $E(n^2)$ in an obvious way. $GL(n)$, the set of invertible $n \times n$ matrices, is a group that inherits from its embedding in $E(n^2)$ the structure of a manifold. In fact, $GL(n)$ so-considered is a Lie group. And so are all of its most famous subgroups, such as the group of matrices of determinant one, the group of orthogonal matrices, the group of symplectic matrices, etc. Very often, when a Lie group is in play it is one of these groups, or one of the subgroups of the group of invertible $n \times n$ complex matrices.
- So far we have focussed on compact groups. In the case of a real matrix group G , compactness means that G is closed and bounded as a subset of $E(n^2)$. The prototypical compact Lie groups are the group of orthogonal matrices (in the real case) and the group of unitary matrices (in the complex case). Here is a definition a principal bundle that works for the non-compact as well as the compact case: a *principal G -bundle* is a manifold P equipped with a smooth, free, proper G -action.

- An action of G on P is *proper* if: for sequences $\{x_n\}$ in P and $\{g_n\}$ in G , the convergence of $\{x_n\}$ and $\{g_n \cdot x_n\}$ in X implies that $\{g_n\}$ has a convergent subsequence in G . Any action of a compact group is proper.
- But what does all of this have to do with the usual sort of definition of a principal bundle? Are these things even fibre bundles in the ordinary sense?
 - Whenever we have a group G acting on a set X , we will want to talk about orbits: the G -orbit of a point x in X is the set $O_x = \{y = g \cdot x : g \in G\}$. We denote by M the set $\{O_x : x \in P\}$ of orbits of the action of G on P . We denote by π the map that sends a point x in X to the corresponding orbit O_x in M .
 - The technical conditions in our definition of a principal bundle ensure everything in sight is smooth and well-behaved. In particular, if P is a principal G -bundle, then: (i) for any x in P , the orbit O_x is a submanifold of P diffeomorphic to G ; (ii) the space of orbits M inherits from P a manifold structure relative to which the map π is smooth; and (iii) for any open set U in M , $\pi^{-1}(U)$ is isomorphic, as a manifold, to $U \times G$. We call M the *base-space* of P .
 - So a principal G -bundle P is a fibre bundle with typical fibre G and with base space M that carries a nice action of G .

Principal fibre bundles in Yang-Mills theories.

- *What is a connection?* Let P be a principal G -bundle. Since P is a manifold, at each point x in P we have the tangent space $T_x P$. Since P is a principal bundle, at each x there is a distinguished subspace V_x of $T_x P$ consisting of those vectors tangent to the submanifold $O_x \subset P$. A *connection on P* is a gadget that selects at each x in P a subspace H_x complementary to V_x (so that each v in $T_x P$ can be written in a unique way as a sum of a vector in H_x and a vector in V_x), in a fashion that is smooth and G -invariant (i.e., if $y = g \cdot x$ then H_y is the image under g of H_x). We can think of a connection on P as adding a bit of geometric structure to our principal bundle (think of the connection as telling us that vectors in H_x are orthogonal to vectors in V_x).
- *What is a connection one-form?* We are used to encoding geometrical structure on manifolds via tensors. We can do that in the present case. A *connection one-form* is a certain sort of (vector-valued) one-form on P . There is a bijective correspondence between connections on P and connection one-forms on P .
- *Bundles and connections in Yang–Mills theory?* As usually formulated, general relativity is a theory whose field is a metric tensor on a manifold representing spacetime; in the same way, as usually formulated, a Yang-Mills theory is a theory whose field is a connection one-form on a principal bundle whose base-space M represents spacetime. In both cases, one is able to ignore the full technical apparatus if one is concerned with special solutions or physical situations; but if one wants to talk about the theory as a whole, one has to face the music (or introduce some equally daunting non-standard formulation of the theory).
- *Why on Earth do we do all this?* Long before they had heard of principal bundles and connections, physicists wrote down field theories that generalized Maxwell's

theory by allowing internal degrees of freedom (isospin etc). This involved writing down equations governing horrible index-laden quantities describing fields on spacetime that obeyed strange transformation laws.

- So far connections are forms on P , a manifold several dimensions larger than spacetime. How do we get quantities defined on spacetime?
 - We proceed as follows. Let U be an open set in the spacetime manifold M . A *local section* over U is a smooth map $s:U\rightarrow P$ such that for any x in U , $\pi(s(x))=x$. The image of such an s will be a submanifold of P diffeomorphic to U . Whenever we have a diffeomorphism d between manifolds, we can use it to pull back tensors defined on one manifold to the other. Let us do that here: take a connection one-form on P ; restrict it to $s(U)\subset P$, then use s to pull back this (vector-valued) one-form to U . The result is a horrible index-laden quantity on spacetime (some indices are spacetime indices, some correspond to the vector space in which the connection one-form takes its values). This looks just like the sort of field that turned up when physicists generalized Maxwell's theory.
 - What happens if we use some other local section $s^*: U\rightarrow P$ to pull back the connection one-form to spacetime? Then we get a different, horrible index-laden quantity on U that arises from the first one we constructed by applying the sort of transformation laws that turned up when physicists generalized Maxwell's theory.
- So it looks like a connection one-form on a principal bundle is the natural intrinsic way to characterize what lies behind all these index-laden quantities, in much the same way that a metric tensor is the natural intrinsic quantity standing behind the index-laden formulations of GR.
- *Symmetries*. Subtleties aside, the symmetry group of general relativity is the group of diffeomorphisms of the spacetime manifold: if d is such a diffeomorphism and g is a solution, then so is d^*g . The symmetry group of a Yang–Mills theory is similarly huge. A *vertical bundle automorphism* is a diffeomorphism from P to itself that respects the G -action on P . This is a huge subgroup of the group of diffeomorphisms of P . Vertical bundle automorphisms map connection one-forms satisfying the Yang–Mills equations to connection one-forms satisfying the Yang–Mills equations.
- *What about matter fields?* Let V be a manifold on which G acts smoothly. A configuration of a matter field for a Yang–Mills theory is represented by a function $\Phi:P\rightarrow V$ such that for any x in P and any g in G , $g^{-1}\cdot\Phi(x)=\Phi(g\cdot x)$ [note that the dots on the two sides of this equation correspond to G -actions on *different* spaces].
 - Usually, one takes V to be a vector space and takes the action of G on V to be linear.
 - Often, configurations of matter fields are taken to sections of a vector bundle with typical fibre V associated with P . The present way of proceeding is equivalent, but saves us introducing further bundles.