# Aharonov-Bohm in 3D

Alexander Afriat

February 20, 2025

**Abstract**

The old vector calculus of `div`, `grad`, `curl` misleadingly distinguishes between two theorems that are only apparently different: Stokes ($\mathbf{St}\,\Theta_2$)—the one needed for the Aharonov-Bohm effect ($\mathbf{AB}_2$)—and Gauss ($\mathbf{St}\,\Theta_3$). With differential forms, the two vector calculus theorems get brought together into a single abstract 'Stokes theorem' ($\mathbf{St}\,\Theta_N$) which holds in two or more dimensions and preserves Stokes's name alone. Since the theorem provides the basic mathematical structure for Aharonov-Bohm, the effect itself ($\mathbf{AB}_2$) must be just as capable of abstract generalisation ($\mathbf{AB}_N$) and three-dimensional reformulation ($\mathbf{AB}_3$) as the theorem; here I propose to understand $\mathbf{AB}_3$ in terms of Newtonian gravity—or equivalently, in terms of electrostatics. By disentangling essential (boundary, hole, flux, source, primitive *etc.*) from accidental (solenoid, screen *etc.*) features, this sheds light on the effect and its interpretations.

## 1 Introduction

We can begin with a string $\alpha_1 \to \alpha_2 \to \alpha_3$ of three issues, each one leading to the next.

$\alpha_1$: By what law does magnetism deflect a charge $e$? What's surprising is the dependance on the size of the charge: one law (referring explicitly to the magnetic field $\mathbf{S}_2 = d\mathbf{P}_2$) for large charges, another (which ignores the magnetic field, referring only to one or other of its primitives $\mathbf{P}_2, \mathbf{P}'_2 \in [\mathbf{P}_2] := d^{-1}d\mathbf{P}_2$) for small ones.

The derivative relating the two laws, for large and small charges, leads directly to the second issue, $\alpha_2$: the troublesome freedom produced by the kernel of a differential operator. One encounters it everywhere, but here I'll write[1] $[\mathbf{P}] = d^{-1}d\mathbf{P}$, usually with a subscript indicating the degree of the differential form, as in $\mathbf{P}_2$ or $\mathbf{P}_3$. So not only do we have very different magnetic deflection laws for large and small charges, which is awkward enough already; but the law for small charges is complicated by the ambiguity $[\mathbf{P}_2] = d^{-1}d\mathbf{P}_2$ characteristic of all primitives (which is overcome by the differentiation yielding the law for large charges).

Which brings us to the third issue, $\alpha_3$: the elimination or even just attenuation of the 'differential' freedom $\alpha_2$, which may not always be as straightforward as one wants. It would be nice if the matter were entirely black & white: *here* we can pick out a single element of the equivalence class $[\mathbf{P}] = d^{-1}d\mathbf{P}$, *there* we're stuck with the whole class, and that's all there is to it. Needless to say, it isn't that simple; leaving out the black and the white as unrealistically binary, we're left with the various shades of grey in between—just about everywhere, in all cases. What may look, here and there,

---

[1] Where $\mathbf{P}$ is a differential form and $d$ the exterior derivative, see page 5 below.

like black or white are no more than deceptively extreme shades of grey enclosing a whole spectrum of complications.

The experiment (§2.2) proposed by Aharonov & Bohm[2] is an attempt to answer $\alpha_1$: the transition from Hamiltonian mechanics to quantum theory seems to *impose* the 'momentum' coupling law

$$p \mapsto p + \mathbf{P}_2$$

(or rather (9) below), hence ruling out an interaction, such as

$$\omega \mapsto \omega + \mathbf{S}_2$$

(or rather (7) below), referring directly to the magnetic field $\mathbf{S}_2$; where $p$ is momentum and $\omega$ the symplectic two-form (§2.1). But how can quantisation exclude a law which was *equivalent* to (9) in Hamiltonian mechanics? How can a magnetic coupling law ignore the magnetic field itself? This is surprising enough to deserve the experimental verification proposed by Aharonov & Bohm (1959). But the momentum law (9) is just background, to get us started, and not the main subject of this paper, which is about foundational questions raised by the Aharonov-Bohm (**AB**) experiment devised to test it; the 'coupling issue' $\alpha_1$ essentially serves to bring up the 'differential issue': how can the charge interact with the magnetic *potential* $\mathbf{P}_2$ rather than its less ambiguous derivative

$$(1) \qquad\qquad d(\mathbf{P}_2 + d\lambda_2) = \mathsf{B} =: \mathbf{S}_2,$$

which is spared the gauge freedom $\alpha_2$ that so notoriously vitiates $\mathbf{P}_2$, casting doubt on its very existence? The magnetic field $\mathbf{S}_2$ and equivalence class

$$(2) \qquad\qquad [\mathbf{P}_2] = d^{-1}d(\mathbf{P}_2 + d\lambda_2)$$

are measurable, but how about[3] $\mathbf{P}_2$ itself? Again, this leads to $\alpha_3$: is measurability really as straightforwardly binary, black & white, as can be hoped?—this *is* measurable, that *isn't*, and that's all there is to it? If it were that straightforward, $\alpha_2$ might be just as straightforward: *here* the differential freedom is altogether removed by measurement, *there* it isn't, and no more need be said. But even if $\alpha_3$ goes with $\alpha_2$, this is clearly no place to propose a general theory of measurement that sorts out $\alpha_3$ once and for all; all I can do here is point out that the reduction of the freedom $\alpha_2$ can be problematic, and subject to the confusing shades of grey on which, as we'll soon see, the 'Stokesian' isomorphism $\mathbf{AB}_2 \sim \mathbf{AB}_3$ captured in the **Table** below ultimately rests.

The awkward 'measurability issue' $\alpha_3$ takes us to a magnetic application

$$\mathbf{F}_2 = \oint_{[\partial \mathbf{L}_2]} [\mathbf{P}_2] = \iint_{[\mathbf{L}_2]} \mathbf{S}_2$$

of Stokes's theorem $\mathbf{St}\,\Theta_2$, which suggests that the problematic field $\mathbf{P}_2$ radiated from the 'source'[4] $\mathbf{S}_2$ somehow reaches the boundary $\mathbf{B}_2 := \partial \mathbf{L}_2$ *without gain or loss* (having meanwhile crossed an annular region defined in §2.2 in which $d\mathbf{P}_2$ vanishes and hence $\mathbf{P}_2$ is *neither produced nor destroyed*); where the equivalence classes $[\partial \mathbf{L}_2]$,

---

[2] See Aharonov & Bohm (1959), but also Franz (1939, 1940, 1965), Ehrenberg & Siday (1949), Olariu & Popescu (1985), Hiley (2013).

[3] In §3.2.2 we'll see that Aharonov & Bohm (1959) propose a way—whose feasibility is not considered here—of measuring $\mathbf{P}_2$.

[4] What I mean by "source" will be clarified in this Introduction, and in §3.1 below.

[$\mathbf{P}_2$] indicate that the boundary $\mathbf{B}_2$ and radiation $\mathbf{P}_2$ can be independently subjected to dual (§2.3.5) deformations without affecting the flux $\mathbf{F}_2$. Details below in §2.2; in §2.3 we'll see that the riches [$\mathbf{P}_2$] can be found embarrassing enough to be altogether *renounced* ($\mathbf{AB}_2$-1), or even better *replaced* ($\mathbf{AB}_2$-3) with Jesuitical subtlety by the homotopy class $\mathbf{H}_2 := [\mathbf{B}_2]$, whose *equivalent* riches—penury is best avoided—are ingeniously disguised to make them less troubling. The special case $\mathbf{St}\,\Theta_2$ then leads to the $N$-dimensional Stokes theorem $\mathbf{St}\,\Theta_N$

$$\mathbf{F}_N = \int \cdots \int_{\partial \mathbf{L}_N} \mathbf{P}_N = \int \cdots \iint_{\mathbf{L}_N} d\mathbf{P}_N,$$

which in turn brings us to this paper's title: however unsettling one may find (the usual two-dimensional) $\mathbf{AB}_2$, however embarrassing one may find the riches [$\mathbf{P}_2$], one can take comfort in the reassuring Stokesian isomorphism $\mathbf{AB}_2 \sim \mathbf{AB}_3$, which suggests that $\mathbf{AB}_2$ may not be too disturbing after all—having as it does so much in common with mere electrostatics ($\mathbf{ES}$) and Newtonian gravity ($\mathbf{NG}$), sound (equivalent) theories that could hardly be less teratological. Indeed in §3.1 we'll see how $\mathbf{AB}$ can be generalised from two ($\mathbf{AB}_2$) to $N$ dimensions ($\mathbf{AB}_N$), and in particular three, $\mathbf{AB}_3$, which is well exemplified by $\mathbf{ES}$ or $\mathbf{NG}$. The Stokesian isomorphism $\mathbf{AB}_2 \sim \mathbf{AB}_3$ is at worst slightly tempered, but not seriously threatened, by the four possible 'obstacles to isomorphism' $O_1$-$O_4$ listed on page 4 below.

The strategy is reminiscent of a *relative consistency proof*,[5] in which a theory (one isn't sure about) is shown to be no less consistent than another (one's used to); by isomorphism, $\mathbf{AB}_2$ will accordingly be shown to be *no more teratological than* $\mathbf{NG}/\mathbf{ES}$. *Holonomic* (§2.3.3) and *topological* (§2.3.4) can be viewed as 'values of the variable' *teratological*.

So the isomorphism strategy rests on the $N$-dimensional Stokes theorem $\mathbf{St}\,\Theta_N$ with differential forms—$\mathbf{St}\,\Theta_2$ (2D) and $\mathbf{St}\,\Theta_3$ (3D) being no more than two different 'values' of the same $N$-dimensional theorem $\mathbf{St}\,\Theta_N$. The low-dimensional versions $\mathbf{St}\,\Theta_2$ and $\mathbf{St}\,\Theta_3$ both correspond to theorems in the old vector calculus of `div`, `grad`, `curl`:

$\mathbf{St}\,\Theta_2$ to the Kelvin-Stokes '`curl` theorem' involving circulation around a loop

$\mathbf{St}\,\Theta_3$ to the Gauss-Green '`div` theorem' involving flux through a closed surface.

In the standard two-dimensional Aharonov-Bohm scheme (§2.2), where the two parts of the split wavefunction are made to interfere on a screen after encircling the solenoid, the Kelvin-Stokes theorem $\mathbf{St}\,\Theta_2$ essentially relates the 'turbulence' $\mathbf{P}_2$ leaving the 'source' $\mathbf{S}_2$ to the turbulence eventually caught, neither diminished nor augmented, by an encircling loop $\mathbf{B}_2$—since the conservation condition $d\mathbf{P}_2 = 0$ holding throughout the intervening annulus $\mathbf{D}_2$ rules out the creation or destruction of $\mathbf{P}_2$. But since $\mathbf{St}\,\Theta_2$ and $\mathbf{St}\,\Theta_3$ are essentially the same theorem ($\mathbf{St}\,\Theta_N$), surely there has to be a 3D version (§3.1) of $\mathbf{AB}$; that, in a nutshell, is the point. To distinguish, I've included a 'dimensional' subscript: so $\mathbf{AB}_2$ is the usual two-dimensional Aharonov-Bohm (with the wavefunction encircling the solenoid), and $\mathbf{AB}_3$ the 3D version (which is captured well enough by $\mathbf{NG}$ and $\mathbf{ES}$—but without a wavefunction or a screen, or even a solenoid).

$\mathbf{AB}_3$ is intended to be no more than an *analogy* and by no means an *exact rendering* of $\mathbf{AB}_2$ in three dimensions—which would in any case be impossible, even

---

[5]See Poincaré (1917) pages 56ff., Weyl (2009) §4.

pointless. Indeed one has to distinguish between the formalism and (an unduly literal understanding of) the physics: the overly literal physics can be misleading, whereas the analogy works well at a more abstract, formal level. The Ampèrian `analogy` proposed at the end of §2.3.4 resembles $\mathbf{AB}_2$ even more, but that may or may not make it better than $\mathbf{AB}_3$; in any case we have three alternative ways of thinking about $\mathbf{AB}_2$: Ampère's law (which I've chosen to exploit very little), gravity, electrostatics.

Reformulation in three dimensions involves translation of the main entities; some of the correspondences can be surprising: the solenoid's magnetic field $\mathsf{B} =: \mathbf{S}_2$ in $\mathbf{AB}_2$ becomes the mass/charge density $\rho =: \mathbf{S}_3$ in $\mathbf{AB}_3$, the magnetic potential $\mathsf{A} =: \mathbf{P}_2$ ($\mathbf{AB}_2$) becomes the gravitational/electrostatic ($\mathbf{Gr/Es}$) field $E =: \mathbf{P}_3$ ($\mathbf{AB}_3$) and so on. But once a 3D ($\mathbf{NG/ES}$) counterpart is considered alongside the standard 2D version of the effect, it makes sense to add a third scheme $\mathbf{AB}_N$ made up of the abstractions common to—indeed distilled from—the other two. For instance: loop ($\mathbf{AB}_2$), membrane ($\mathbf{AB}_3$), boundary ($\mathbf{AB}_N$); electromagnetic field ($\mathbf{AB}_2$), mass/charge density ($\mathbf{AB}_3$), source ($\mathbf{AB}_N$); electromagnetic potential ($\mathbf{AB}_2$), gravitational/electrostatic field ($\mathbf{AB}_3$), primitive ($\mathbf{AB}_N$).

**Table**

| $\mathbf{AB}_2$ | $\mathbf{AB}_3$ | $\mathbf{AB}_N$ |
|---|---|---|
| magnetic field $\mathbf{S}_2 := \mathsf{B}$ | source density $\mathbf{S}_3 := \rho$ | source $\mathbf{S}$ |
| mag. potential $\mathbf{P}_2 := \mathsf{A}$ | $\mathbf{Gr/Es}$ field $\mathbf{P}_3 := E$ | primitive $\mathbf{P}$ |
| equiv. class $[\mathsf{A}] = d^{-1}\mathsf{B}$ | equiv. class $[E] = d^{-1}\rho$ | equiv. class $[\mathbf{P}] = d^{-1}\mathbf{S}$ |
| region $\mathbf{L}_2$ | region $\mathbf{L}_3$ | region $\mathbf{L}$ |
| loop $\mathbf{B}_2 := \partial\mathbf{L}_2$ | membrane $\mathbf{B}_3 := \partial\mathbf{L}_3$ | boundary $\mathbf{B} := \partial\mathbf{L}$ |
| hoop $\mathbf{H}_2 := [\mathbf{B}_2]$ | homot. class $\mathbf{H}_3 := [\mathbf{B}_3]$ | homot. class $\mathbf{H} := [\mathbf{B}]$ |
| support $\mathbf{I}_2$ of $\mathbf{S}_2$ | support $\mathbf{I}_3$ of $\mathbf{S}_3$ | support $\mathbf{I}$ of $\mathbf{S}$ |
| inner circle $\partial\mathbf{I}_2$ | inner sphere $\partial\mathbf{I}_3$ | inner sphere $\partial\mathbf{I}$ |
| larger disc $\mathbf{J}_2$ | larger ball $\mathbf{J}_3$ | larger ball $\mathbf{J}$ |
| outer circle $\partial\mathbf{J}_2$ | outer sphere $\partial\mathbf{J}_3$ | outer sphere $\partial\mathbf{J}$ |
| annulus $\mathbf{D}_2$ | difference $\mathbf{D}_3$ | difference $\mathbf{D}$ |
| circulation $\mathbf{F}_2$ | flux $\mathbf{F}_3$ | flux $\mathbf{F}$ |
| `curl` | `div` | exterior derivative $d$ |
| magnetic field interpret. | source interpret. | source interpret. |
| potential interpret. | $\mathbf{Gr/Es}$ field interpret. | primitive interpret. |
| holonomy interpret. | membrane interpret. | boundary interpret. |
| topology interpret. | topology interpret. | topology interpret. |

Again, there come to mind four potentially awkward differences between $\mathbf{AB}_2$ and $\mathbf{AB}_3$, that could somewhat temper the isomorphism. The first possible obstacle $O_1$, which takes us to the measurability issue $\alpha_3$, is neither uninteresting nor unsurmountable; the other three seem less interesting or threatening, especially in a mere thought experiment, and aren't worth dwelling on.

$O_1$ The field $\mathbf{P}_3$ is usually taken to be a good deal more measurable than $\mathbf{P}_2$; one gets the impression that whereas a single (Coulombian?) element can be picked out of the class $[\mathbf{P}_3] = d^{-1}d\mathbf{P}_3$, there's no way of overcoming the freedom $[\mathbf{P}_2] = d^{-1}d\mathbf{P}_2$—one's just stuck with the whole awkward, unwieldy class.

$O_2$ In $\mathbf{AB}_2$, the source $\mathbf{S}_2$ (and with it the flux $\mathbf{F}_2$) can be strengthened differentially, gradually, adiabatically—but through the other spatial dimension (the third

one), which is usually assumed away. Where can a further *spatial* dimension, a fourth one, be found in $\mathbf{AB}_3$?

$O_3$ In $\mathbf{AB}_2$, interference fringes on the screen indicate the flux $\mathbf{F}_2$ through the boundary $\mathbf{B}_2$. But how does one measure the flux $\mathbf{F}_3$ through the two-dimensional boundary $\mathbf{B}_3$?

$O_4$ An appropriate exactness condition $\mathbf{EX}_2$ (analogous to $\mathbf{EX}_3$) seems unavailable.[6]

In $\mathbf{AB}_2$, the trouble boils down to $\alpha_2 \to \alpha_3$: the magnetic field (1) is measurable but the gauge freedom (2) is taken[7] to be insuperable. For the analogy to hold, the corresponding freedom

$$(3) \qquad\qquad [\mathbf{P}_3] = d^{-1}d(\mathbf{P}_3 + d\lambda)$$

in $\mathbf{AB}_3$ would have to be about as hard to overcome. The catch ($O_1$) being that the analog $\mathbf{P}_3$ of $\mathbf{P}_2$ is the gravitational/electrostatic field $E$, which is often considered *trivially measurable*. But in §3.2.1 I'll argue that $\mathbf{P}_3$ in $\mathbf{AB}_3$ isn't *that much* more measurable than $\mathbf{P}_2$ in $\mathbf{AB}_2$ after all; Aharonov & Bohm (1959) even suggest (§3.2.2) that $\mathbf{P}_2$ may be *measurable* ($\alpha_3$). And even if $\mathbf{P}_3$ isn't every bit as unmeasurable as the magnetic potential, what's usually claimed is not that the spherical symmetry of the inverse square law—or Coulomb's law $\mathbf{CL}$—is obtained by experimental verification at each and every point of $\mathbb{R}^3$, which would be tiresome, indeed impossibly so; but rather that it can be derived from Poisson's equation by *controversially*[8] assuming the asymptotic constancy $\mathbf{AC}$ of the potential $\phi$ (which by figuring in Poisson's equation is assumed to exist). Spherical symmetry can perhaps be derived from the *principle of sufficient reason*—which rests, however, on unrealistically strong homogeneity assumptions; what happens to it against an awkwardly inhomogeneous *curved* background? What happens to $\mathbf{CL}$ if the constitutive properties[9] of the medium are allowed to vary in a very general way? The principle can only guarantee spherical symmetry or isotropy against an ideal, perfect background, whose irregularities would otherwise affect the fields.

However one feels about Coulomb's law, spherical symmetry, the principle of sufficient reason, constitutive relations, curvature or the measurability of the electrostatic field (§3.2.1 below), the neglected relationship $\alpha_2 \to \alpha_3$ between gauge freedom and measurability deserves attention and needs to be spelled out—it isn't enough to go on repeating the familiar

slogan: *(only) gauge invariant quantities are real*

(or objective or whatever), which seems to do no more than beg the question, putting the cart before the horse. And it could be—but here I'll only raise the issue, this isn't the place to dwell on it—that we're drawing the wrong distinctions and hence 'classifying' gauge freedoms inappropriately: maybe the magnetic gauge freedom (2) ultimately has more in common with the formally identical gravitational/electrostatic freedom (3) considered below than with, for instance, the coordinate freedom of general relativity.[10]

---

[6] See §3.2.1 below.

[7] *Cf.* §3.2.2 below.

[8] Einstein (1917) is all about such asymptotic assumptions, see footnote 49 below. He finds them so troubling he goes so far as to eliminate spatial infinity, proposing a new topology with only temporal infinities.

[9] Quite generally, these are the susceptibilities expressed by the Hodge duality $*$ between electromagnetism represented by the Faraday two-form $F$ and $G = *F$, which Kottler (1922) calls *Magnetelektrismus*; here we can think of the dielectric tensor $\varepsilon$.

[10] But see §2.3.5 below, where magnetic gauge transformations are literally viewed as diffeomorphisms.

The inverse image $d^{-1}$, which first appeared in (2) above, can be problematic elsewhere, but not here, where it works well and makes sense. By 'expanding' the differential form $\mathbf{P}$ to its gauge equivalence class $[\mathbf{P}] = d^{-1}d\mathbf{P}$, the composition $d^{-1}d$ gives expression to the derivative's (nontrivial) kernel, which is behind all the trouble here ($\alpha_2$). I've put the exact term $d\lambda$ on the right, where it is annihilated by the $d$ before the round parentheses; but it can also be put on the left, so that one sees the equivalence class

$$(4) \qquad\qquad [\mathbf{P} + d\lambda]_\lambda = d^{-1}d\mathbf{P}$$

obtained varying the differential form $\lambda$ which, much like a constant of integration, serves to 'restore the loss caused by $d$'; $\lambda_2$ being a 'function of integration,' $\lambda_3$ a 'one-form of integration' and so on. So each element of (4) corresponds to a different $\lambda$; a particular $\lambda$ returns the original $\mathbf{P}$.

But the 'noninvertibility' or 'nullspace' problem ($\alpha_2$) behind all the gauge freedom is so basic it doesn't even require derivatives, limits, analysis: mere subtraction $x_2 - x_1$ is enough. Applied to two real numbers $x_1, x_2$, the difference operator

$$\Delta = \Delta(\,\cdot\,,\cdot\,) : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$$

yields the difference

$$\delta = \Delta(x_1, x_2) = x_2 - x_1;$$

but the inverse image $\Delta^{-1}\delta$ is the set, with cardinality $\mathbb{R}$, of all intervals of width $\delta$ (or more precisely the set of all pairs $x_1, x_2$ such that $x_2 - x_1 = \delta$). Fixing a particular interval—the original one, for instance—requires a specification (*e.g.* a particular first element $\bar{x}_1$) along the lines of a constant of integration. The composition $\Delta^{-1}\Delta$ maps a specific interval $(x_1, x_2)$ to the set of all intervals with width $\delta = \Delta(x_1, x_2)$.

More on $O_1$ in §3.2; as for $O_2$, it can be turned around to reinforce, rather than weaken, the Stokesian isomorphism $\mathbf{AB}_2 \sim \mathbf{AB}_3$. In $\mathbf{AB}_2$, the magnetic field gets strengthened 'through the additional spatial dimension' as it were (the third one), in the sense that the current climbs up the $z$ axis around which the coil is wound. Even if a two-dimensional 'magnetic field' were somehow conceivable, its modulations wouldn't be without the third dimension. In $\mathbf{AB}_3$, the conservation of charge rules out its augmentation; but much as in $\mathbf{AB}_2$, an additional spatial dimension (to which one can presumably help oneself in a mere thought experiment) would get around the problem by allowing charge to be gradually 'fed into $\mathbb{R}^3$' through that fourth spatial dimension; conservation would only be violated in 3D, not in 4D. Not science fiction, just a sober thought experiment taking advantage of the dimensional possibilities provided by the analogy itself. But if one disapproves of four spatial dimensions, $O_2$ is entirely dealt with by the Ampèrian `analogy` at the end of §2.3.4 (which almost resembles $\mathbf{AB}_2$ *too* much), where the source gets strengthened through the *third* spatial dimension—which one's surely entitled to, even within ordinary 3D physics.

Overcoming the third potential obstacle, $O_3$, will be left up to the ingenuity of the thought-experimenter. If it is objected that measurement of the *single* number $\mathbf{F}_3$ could be so hard as to require *much* ingenuity, the

$$\infty^3 = \infty_x \times \infty_y \times \infty_z$$

numbers $\mathbf{P}_x$, $\mathbf{P}_y$, $\mathbf{P}_z$ (the values of the three functions) considered in §3.2.1 below must be that much worse.

6

Which brings us to $O_4$: Even if an asymptotic condition along the lines of **AC**—which would significantly *constrain* A—could be formulated, I cannot see how one could then finish the job, *fixing* A with an exactness condition (**EX**$_2$) analogous to **EX**$_3$:

$$\mathbf{P}_3 = *d\phi,$$

where the star denotes Hodge duality. Since there's nothing below a zero-form (the potential $\phi$), a two-dimensional version **EX**$_2$ of **EX**$_3$ seems unavailable (§3.2.1).

But all told the isomorphism $\mathbf{AB}_2 \sim \mathbf{AB}_3$ survives broadly intact, only slightly tempered by $O_1$-$O_4$.

In the literature one finds three or even four main interpretations of $\mathbf{AB}_2$; going from two to three dimensions, then even to the abstract scheme $\mathbf{AB}_N$ common to both, sheds light on all four. Though two of the interpretations sketched in §2.3 rely on the unmeasurability of the primitive $\mathbf{P}$, the topology interpretation $\mathbf{AB}$-4 doesn't seem to at all—which would confirm that $\mathbf{AB}_2$ is *no more topological than* $\mathbf{AB}_3$. And *no more topological than* $\mathbf{AB}_3$ could be taken to mean *uninterestingly topological* or even *hardly topological at all*, since one seldom bothers even to notice the trivially topological character of Newtonian gravity or electrostatics.

But this paper is more about the abstract mathematical structure and fundamental nature of the Aharonov-Bohm effect than about particular interpretations (in two dimensions), whose zealous partisans may see no more than hasty caricatures in my deliberately summary accounts. Even if interpretations are clearly worth mentioning, I have not pursued exegetical accuracy, this is by no means a full review or survey of them, their details can be found elsewhere; I have sought rather to identify some of them concisely (without going to the trouble of concealing my sympathies) than to plead. I imagine far too many have been proposed over the decades for a complete enumeration to be even attempted.

The trouble with the standard two-dimensional treatment $\mathbf{AB}_2$ is that on its own it confusingly entangles essential and accidental features, thus vitiating and perplexing the interpretations. The abstract $\mathbf{AB}_N$ column of the above **Table** represents something like 'essential features disentangled from accidental features.' A logician may even see a *theory* in the third column alongside two of its *models*, represented in the first two columns. The wavefunction and screen are rather inessential features of the two-dimensional case, as are the *peculiarities* of the source (the solenoid); *a* source, its primitive(s) and enclosing boundary are essential. What's left in the third column is clearly rather bare and mathematical; but at least one's attention is directed to the fundamental structures and logic, without being distracted by confusing details of doubtful relevance.

Abstraction, certainly in the third column of the **Table**, is taken to a level going well beyond the idealisations thoroughly considered in Shech (2015, 2018a) and Earman (2017). The solenoid doesn't even bother being infinitely long, it is just replaced by an abstract source; boundaries are perfectly impenetrable and so on. Interesting, fundamental aspects of the Aharonov-Bohm effect and debate are situated at this rather high level of abstraction. Other aspects—none the less interesting and fundamental—are of course situated at the lower level(s) of abstraction examined in Shech (2015, 2018a) and Earman (2017). We're all familiar with the 'horizontal' choice of subject or discipline: electromagnetism or mechanics or botany or history; but an investigation or analysis also requires the 'vertical' choice of a level of abstraction. Both choices, horizontal and vertical, involve inclusion, omission, emphasis and neglect. Aspects considered at one level of abstraction may not figure at another, that's just part of the

choice.

Accustomed to abstractions, the geometer can find the move from

$$\mathbf{AB}_2\&\mathbf{St}\,\Theta_2 \text{ to } \mathbf{AB}_3\&\mathbf{St}\,\Theta_3 \text{ then } \mathbf{AB}_N\&\mathbf{St}\,\Theta_N$$

rather trivial (without, admittedly, being too far off the mark—but the opposite danger, of seeing no analogy at all, is surprisingly the real problem). That would, however, amount to pushing abstraction a bit too far for our purposes, leaving *too much* out: even if some details are irrelevant, confusing, misleading, others are significant and indeed welcome; though $\mathbf{St}\,\Theta_2$ is just $\mathbf{St}\,\Theta_N$ with $N=2$, there's more to $\mathbf{AB}_2$ and $\mathbf{AB}_3$ than just $\mathbf{St}\,\Theta_N$. The middling level of abstraction I choose is lower than 'just $\mathbf{St}\,\Theta_N$ on its own,' while remaining above that of $\mathbf{AB}_2$ in all its gory detail.

So the three-dimensional analogy is intended to shed light on a number of foundational and philosophical issues: non-locality, the relation $\alpha_2 \to \alpha_3$ between empirical accessibility and gauge freedom, the physical status of mathematical boundaries, the role of topology and so forth. It is useful to see, for instance, that $\mathbf{AB}_2$ is *no more topological* than gravity or electrostatics; and that in three dimensions, hoops of loops become equivalence classes of two-dimensional membranes. Many of the philosophical implications may well deserve to be spelled out explicitly; but the main point of the analogy is more abstract: to change one's *general* perspective, and understanding of the effect, which could in principle have countless philosophical implications, more than can be enumerated here.

Wallace (2014, page 15) chooses[11] to view $\mathbf{AB}_2$ as an effect involving a single "jointly described" quantomagnetic entity:

> But how are we to think about this "jointly described" entity? We know that it can be characterised entirely by the magnitude of $\psi$ (a scalar field) and by its covariant derivative (a vector field, or more precisely a one-form field). It is important to remember that these are conceptually and mathematically very different entities. A scalar field, mathematically, is just an assignment of a real number to every point of space, and can easily enough be though of as ascribing properties to *points* of space. A one-form field is not so simple and cannot be so represented: to speak loosely, it is more like an assignation of properties to infinitesimally small differences between points of space.

The interaction between the wavefunction and the potential $\mathbf{P}_2$ is hardly uninteresting; without $\mathbf{P}_2$ one would know neither how to differentiate nor how to compare neighbouring values of the wavefunction (nor even distant values for that matter, which presuppose step-by-step infinitesimal comparison along a path); the connection, the covariant derivative (essentially $\mathbf{P}_2$) provides a 'horizontal' notion of constancy to which differences can then be referred. But here I've chosen to concentrate on magnetism and hence to view the wavefunction, its phase and the interference pattern as *entities that test* $\mathbf{P}_2$, or rather the whole class $[\mathbf{P}_2] = d^{-1}d\mathbf{P}_2$. As Wallace points out, I'm not alone in this; the Aharonov-Bohm experiment ("Significance of electromagnetic

---

[11]Page 2: "In this paper I argue that much of this debate rests upon a mistake: that of considering the A-field in isolation rather than in conjunction with the $\psi$-field. [. . .] In section 4 I show in a different way how this apparent nonlocality arises in the study of A alone and how it is blocked when we allow for A and $\psi$ jointly." And page 8: "The A-B effect arises because of certain features of the mathematical theory of a complex scalar field $\psi$ coupled to a real vector field A. It is therefore in hindsight a little odd that the literature on the A-B effect has been almost wholly concerned with the A field and hardly at all with the $\psi$ field."

potentials in the quantum theory") is often, perhaps even typically viewed as an *(electro)magnetic experiment concerning the (electro)magnetic potential*—rather than as a *quantum-mechanical experiment meant to test the phase of the wavefunction.*

To avoid the misunderstandings to which my proposal somehow seems to lend itself, it's worth emphasising what I'm *not* doing: however hard one may try to see the (rightly neglected) *vertical direction* somewhere in my title, I'm by no means extending the standard two-dimensional treatment $\mathbf{AB}_2$ to three dimensions *by merely rehabilitating that vertical direction*, 'up the solenoid.' What I'm proposing is not slightly but *completely* different. The $\mathbf{AB}_2$ effect is *rightly* discussed in the literature as two-dimensional—one thinks of a horizontal section, an $xy$ plane.[12] Most of the relevant features vary relatively little along the vertical direction $z$; even if some (such as the wavefunction) can vary more, the emphasis on two dimensions *remains entirely legitimate* as far as I can tell—and indeed the Stokes theorem that's used is the two-dimensional version $\mathbf{St}\,\Theta_2$ involving the `curl` and circulation around a loop. So

$\mathbf{AB}_2$ is just the standard formulation of the effect

$\mathbf{AB}_3$ involves $\mathbf{NG/ES}$ (rather than the usual solenoid and screen *etc.*).

Again, the vertical direction in $\mathbf{AB}_2$—up the axis of the solenoid—is of little relevance here (despite fleeting appearances above, and towards the end of §2.2), and is not being rescued or rehabilitated or reconsidered.

Mathematicians can prefer to represent the structures at issue here in terms of homology & cohomology. But since differential forms—with the associated notions of closed, exact, perforated—are (not only needed but) enough to capture everything relevant, I have chosen to avoid the demanding abstractions of de Rham theory.

Summing up, my strategy relies on

- isomorphism $\mathbf{AB}_2 \sim \mathbf{AB}_3$

    - $\mathbf{St}\,\Theta_N$ with $N = 2,3$
    - formally identical freedoms $[\mathbf{P}_2] = d^{-1}d\mathbf{P}_2$ and $[\mathbf{P}_3] = d^{-1}d\mathbf{P}_3$
    - $O_1$-$O_4$ not threatening enough to get in the way of $\mathbf{AB}_2 \sim \mathbf{AB}_3$

- relatively straightforward, unproblematic 'normality' of $\mathbf{NG/ES}$ ($\mathbf{AB}_3$).

In §2.2 I briefly describe relevant features of the Aharonov-Bohm effect, in §2.3 the four 2D interpretations. In §3.1 I propose the gravitational/electrostatic analogy $\mathbf{AB}_3$, in §3.3 the 3D versions of the four interpretations; having argued in §3.2 that the gravitational/electrostatic field $\mathbf{P}_3$ may not be *that much* more measurable than the electromagnetic potential $\mathbf{P}_2$ after all.

# 2   Two dimensions

## 2.1   The coupling law ($\alpha_1$)

We shall show that, contrary to the conclusions of classical mechanics, there exist effects of potentials on charged particles, even in the region where all the fields (and therefore all the forces on the particles) vanish.[13]

---

[12]The nature of an $n$-form depends on the size of the environment; so a two-form $\alpha_3 \in \bigwedge^2 \mathbb{R}^3$ corresponds to a vector in $\mathbb{R}^3$ but to a density $\alpha_2 \in \bigwedge^2 \mathbb{R}^2$ on the plane. See §2.2 below.

[13]Aharonov & Bohm (1959) page 485 (Abstract)

> The essential result of the previous discussion is that in quantum theory, an electron (for example) can be influenced by the potentials even if all the field regions are excluded from it.[14]

Even if the Aharonov-Bohm effect $\mathbf{AB}_2$ involves (electro)magnetism interacting with a *wavefuction* $\Psi$, we should first consider a *classical* charge. In four-dimensional spacetime $\mathcal{M}$ we have the Lorentz force[15]

$$F^\flat(V) := F(\,\cdot\,, V)$$

perpendicular to the four-velocity $V$ of the (unit) charge $e$, where the Faraday two-form

$$F = F(\,\cdot\,,\cdot\,) \in \bigwedge{}^2 \mathcal{M}$$

represents the electromagnetic field. Though the Aharonov & Bohm (1959) treatment is relativistic and four-dimensional ("$A_\mu(x)$" *etc*.), relativity is of limited relevance here, we might as well confine ourselves to mere magnetism; in three dimensions we have the Lorentz term $v \times \mathsf{B}$, where $v$ is the three-velocity. If the quantum law were similar—referring to $\mathsf{B}$ rather than its primitive $\mathsf{A}$—we'd be spared all the trouble due to the problematic gauge freedom (2), or

$$(5) \qquad\qquad \mathbf{P}_2 \mapsto \bar{\mathbf{P}}_2 := \mathbf{P}_2 + d\lambda_2,$$

where $\lambda_2$ is a zero-form. So how does the transition from classical to quantum mechanics somehow yield an magnetic coupling law involving $\mathsf{A}$ *rather than* its less ambiguous derivative $\mathsf{B}$?

Since quantisation is usually 'Hamiltonian,' we can start with *classical* Hamiltonian mechanics. The Hamiltonian vector field

$$X_{\mathscr{H}} = \omega^\sharp(d\mathscr{H})$$

on phase space $\Gamma$ is obtained by 'symplectically' converting the differential

$$d\mathscr{H} \in \bigwedge{}^1 \Gamma$$

of the Hamiltonian $\mathscr{H} : \Gamma \to \mathbb{R}$, where the phase space can (with no relevant loss of generality) be thought of as the $2N$-dimensional cotangent bundle[16]

$$(6) \qquad\qquad \Gamma = T^*\mathcal{Q}$$

of the $N$-dimensional configuration space $\mathcal{Q}$ (covered by the first $N$ coordinates $q_1, \ldots, q_N$); and $\omega^\sharp(\,\cdot\,)$ is obtained from the inverse $\omega^{-1}(\,\cdot\,,\cdot\,)$ of the *symplectic two-form* $\omega(\,\cdot\,,\cdot\,)$

---

[14]Aharonov & Bohm (1959) page 490

[15]*Cf.* Aharonov & Bohm (1959) page 490: "The Lorentz force [...] does not appear anywhere in the fundamental theory, but appears only as an approximation holding in the classical limit."

[16]Having written (6) I should point out that—especially taken together—the notations $T^*\mathcal{Q}$ and $T\mathcal{Q}$ for the Hamiltonian and Lagrangian bundles are somewhat misleading. To begin with, the only manifold really needed in symplectic geometry (Hamiltonian mechanics) is $\Gamma$, which corresponds to all $2N$ coordinates $q_1, \ldots, q_N, p_1, \ldots, p_N$; coordinate subsets (such as $q_1, \ldots, q_N$) may or may not satisfy the integrability condition of 'corresponding to manifolds.' But let's assume that in both kinds of mechanics there's a manifold $\mathcal{Q}$ for the position coordinates $q_1, \ldots, q_N$. In Lagrangian mechanics, the fundamental space $\mathcal{Q}$ is fixed and subjected to its own diffeomorphisms (which then get extended to $T\mathcal{Q}$ by differentiation); whereas in Hamiltonian mechanics the fundamental space $\Gamma$ is subjected to a (symplectic) class of diffeomorphisms so broad that $\mathcal{Q}$ gets displaced within $\Gamma$, allowing the possibility of $\Gamma = T^*\mathcal{Q} = T^*\mathcal{Q}'$ for very different $\mathcal{Q}$, $\mathcal{Q}'$.

by partial evaluation. The operation $d\mathscr{H} \mapsto X_{\mathscr{H}}$ is something like index lowering in tensor calculus.

Magnetism, which has yet to appear, can do so in (at least) a couple of equivalent ways.[17] The most intrinsic and fundamental is the addition

$$\omega \mapsto \bar{\omega} := \omega + \mathsf{B} \tag{7}$$

of the magnetic two-form $\mathsf{B}$ to the symplectic background $\omega$, yielding the magnetic symplectic two-form $\bar{\omega}$.[18] Still no sign of $\mathsf{A}$ though. Since $\omega = d\theta$ is exact, where $\theta$ is the *canonical one-form*, we can write[19]

$$\bar{\omega} = d\theta + d\mathsf{A} = d(\theta + \mathsf{A}) = d\bar{\theta}, \tag{8}$$

with the magnetic substitution $\theta \mapsto \bar{\theta} = \theta + \mathsf{A}$, where $d\mathsf{A} = \mathsf{B}$.

Even if $\mathsf{A}$ has now appeared, there remains the matter of quantisation; like $\bar{\omega}$, $\bar{\theta}$ is symplectic, whereas quantum mechanics isn't—however Hamiltonian it may look[20] in all its 'canonical' garb: momentum $p$ (rather than velocity $\dot{q}$), Hamiltonian *etc*. At any rate, $\mathsf{A}$ will somehow have to find its way into the Hamiltonian. With respect to canonical coordinates

$$(q, p) := (q_1, \ldots, q_N, p_1, \ldots, p_N)$$

the canonical one-form assumes the peculiar form

$$\theta = \sum_{k=1}^{N} p_k dq_k = \sum_{k=1}^{N} (p_k dq_k + 0 dp_k) \in \bigwedge^1 \Gamma,$$

in which the last $N$ terms, in $dp$, are missing—making it look (despite dependence on all $2N$ coordinates) like a one-form on the $N$-dimensional configuration space $\mathcal{Q}$ instead. But since

$$\mathsf{A} = \sum_k \mathsf{A}_k dq_k \in \bigwedge^1 \mathcal{Q},$$

a mere function of position, *really is* defined on configuration space $\mathcal{Q}$, we can write

$$\bar{\theta} = \sum_k (p_k + \mathsf{A}_k) dq_k,$$

and therefore

$$p_k \mapsto \bar{p}_k := p_k + \mathsf{A}_k, \tag{9}$$

---

[17]Guillemin & Sternberg (1984) page 141: "We have introduced the effect of the electromagnetic field by keeping the same Hamiltonian $\mathscr{H}$ but modifying the symplectic structure. In the standard physics literature, there is another procedure, called "minimal coupling," for obtaining the Lorentz equation that keeps the original symplectic structure but modifies the Hamiltonian. Let us pause to show that the two procedures are formally equivalent." See also Singer (2004) page 71.

[18]*Cf.* Aharonov & Bohm (1959) page 490: "This Schrödinger equation is obtained from a canonical formalism, which cannot be expressed in terms of the [electric & magnetic] fields alone, but which also requires the potentials." But the canonical formalism *can* be expressed in terms of the magnetic field $\mathsf{B}$ alone, as we've just seen—the potential isn't needed.

[19]*Cf.* Aharonov & Bohm (1959) page 485: "It is true that in order to obtain a classical canonical formalism, the potentials are needed." In symplectic terms, this may only mean that the symplectic two-form, being exact, can be written $\omega = d\theta$—allowing us to write (8) as well.

[20]*Cf.* Aharonov & Bohm (1959) page 485: "In the quantum mechanics, however, the canonical formalism is necessary, and as a result, the potentials cannot be eliminated from the basic equations."

which is the rule we're after; in three dimensions[21]

$$p \mapsto \bar{p} = p + \mathsf{A}$$
$$(p_x, p_y, p_z) \mapsto (\bar{p}_x, \bar{p}_y, \bar{p}_z) = (p_x + \mathsf{A}_x, p_y + \mathsf{A}_y, p_z + \mathsf{A}_z).$$

Since the non-magnetic Hamiltonian is a function $(q, p) \mapsto \mathscr{H}(q, p)$ of the non-magnetic canonical coordinates, the magnetic Hamiltonian will be the *same* function $\mathscr{H}(q, \bar{p})$ of the magnetic canonical[22] coordinates $(q, \bar{p})$.

Which leaves quantisation. If the $2N$ coordinates $(q, \bar{p})$ are now appropriately interpreted as quantum magnetic position & momentum operators on Hilbert space, the resulting $\mathscr{H}$ will be the quantum magnetic Hamiltonian; since the abstract non-magnetic momentum operator $\underline{p}$ on Hilbert space becomes the derivative operator $i\partial$ in position space, the abstract magnetic momentum operator $\underline{\bar{p}}$ becomes the covariant derivative $i\partial + \mathsf{A}$ in position space:

| Status | Canonical | Abstract Hilbert space | Position space |
|---|---|---|---|
| non-magnetic | $p$ | $\underline{p}$ | $i\partial = U\underline{p}U^\dagger$ |
| magnetic | $\bar{p} := p + \mathsf{A}$ | $\underline{\bar{p}} := \underline{p} + \underline{\mathsf{A}}$ | $i\partial + \mathsf{A} = U\underline{\bar{p}}U^\dagger$ |

where $U$ is an appropriate unitary operator with adjoint $U^\dagger$. The quantum peculiarity, which everything turns on, is that while the *classical* magnetic coupling could be expressed in terms of $\mathsf{A}$ *or* its derivative[23] $\mathsf{B}$, quantisation yields a more stringent rule, which surprisingly *imposes* the potential $\mathsf{A}$ rather than $\mathsf{B}$—the very possibility of referring the coupling to $\mathsf{B}$ having been somehow discarded, left behind by quantisation.[24] Aharonov & Bohm (1959) interestingly see this as the basis for their whole agenda, and blame the exclusion of $\mathsf{B}$ on the *genuinely* (as opposed to *ostensibly*, even *deceptively*) canonical character they surprisingly attribute to quantum theory. "Canonical" and "symplectic" are practically synonymous—or perhaps "symplectic" is the intrinsic, geometrical notion whereas "canonical" concerns coordinates—but in any case, whatever nuance may separate the two terms, the quantum formalism isn't strictly speaking

---

[21]Born (1925, page 238) shows how (9) can be derived from the Lorentz force law. He (more or less) writes the Lagrange force equation

$$\frac{d}{dt}\frac{\partial M}{\partial \dot{x}} - \frac{\partial M}{\partial x} = -\frac{e}{c}[\mathfrak{v}\mathfrak{H}]_x$$

with $M = \frac{e}{c}\mathfrak{A}\mathfrak{v} = \frac{e}{c}(\mathfrak{A}_x\dot{x} + \mathfrak{A}_y\dot{y} + \mathfrak{A}_z\dot{z})$, where $[\mathfrak{v}\mathfrak{H}]$ is the cross product of the velocity and the magnetic field $\mathfrak{H} = \text{rot } \mathfrak{A}$. From the Lagrangian he derives the Hamiltonian, and (9).

[22]It is worth pointing out that the non-magnetic coordinates $(q, p)$ are canonical with respect to the non-magnetic symplectic form $\omega$ and the corresponding (non-magnetic) Poisson brackets $\{\cdot, \cdot\}_\omega$, whereas the magnetic coordinates are canonical with respect to the magnetic symplectic form $\bar{\omega}$ and the corresponding Poisson brackets $\{\cdot, \cdot\}_{\bar{\omega}}$.

[23]Aharonov & Bohm (1959) page 485: "It is true that in order to obtain a classical canonical formalism, the potentials are needed. Nevertheless, the fundamental equations of motion can always be expressed in terms of the fields alone."

[24]*Cf.* Vaidman (2012, page 040101-3), who avoids the problems—measurability, freedom *etc.*—of potentials by avoiding potentials altogether, with a radically different approach to quantisation: "One might wonder why, instead of performing exact calculations in the framework of quantum mechanics, I consider particles and cylinders pushed by fields in the framework of classical mechanics and then use the correspondence principle to calculate the shifts of the quantum wave packets of particles and cylinders. I have to follow this path because the standard formulation of quantum mechanics, and the Schrödinger equation in particular, are based on potentials. I hope that a general formalism of quantum mechanics based on local fields [F] will be developed. It will provide a solution to the problem of motion of a quantum particle in a force field even if there is no potential from which it can be derived."

canonical, and certainly isn't symplectic. Like $\mathsf{B}$, the symplectic manifold gets discarded by quantisation.

At any rate, the classical interaction (7) could not but vanish where $\mathsf{B}$ did; but since the quantum interaction (9) no longer refers to $\mathsf{B}$, it has no reason to vanish where $\mathsf{B}$ does. Probabilities, statistics are of course affected by (9), and even by the addition to $\mathsf{A}$ of an exact one-form $d\eta$ with components $(\partial_x \eta, \partial_y \eta, \partial_z \eta)$. Statistical invariance only results once

$$\mathsf{A} \mapsto \mathsf{A}' := \mathsf{A} + d\eta$$

is balanced[25] by the phase transformation

$$\Psi \mapsto \Psi' := e^{i\eta}\,\Psi;$$

the same applies to the Lagrangian

$$\mathscr{L}(\Psi, \mathsf{A}) = \mathscr{L}(\Psi', \mathsf{A}') \neq \mathscr{L}(\Psi, \mathsf{A}')$$

yielding the Schrödinger equation as Euler-Lagrange equation.

The new interaction law (9) produced by quantisation is interesting and surprising enough to be worth testing.[26] How? The test has to be an interference experiment, since (9) yields a phase factor

$$\exp i \oint \mathsf{A}.$$

No need to look farther than Young's slits, the most basic, standard interference experiment—duly adapted, essentially by introducing magnetism: a solenoid. Where? It can only be placed 'in the middle,' where it is enclosed by the two beams *but insulated from them*, so that the wavefunction *interacts only with* $\mathsf{A}$*, not with* $\mathsf{B}$.[27] That, in a nutshell, is the Aharonov-Bohm experiment.

The replacement of (7) by (9) is undeniably surprising and experimentally confirmed; but not the real subject of this paper, since I've chosen[28] to take all the quantum mechanics (wavefunction, interference pattern *etc.*) to be *no more than a way of measuring the circulation of* $\mathsf{A}$ *around a loop*; or more abstractly, to be no more than a way of measuring a flux through a boundary. Like much of the standard literature, I've chosen to focus on the nature and very existence of *other* noteworthy peculiarities the experiment devised to test (9) may or may not reveal.

## 2.2 $\mathrm{AB}_2$

I've already sketched the (2D) Aharonov-Bohm experiment above in §2.1, but details still have to be filled in: A wavefunction $\Psi$ is split into two, and these, having enclosed a (deformable) region $\mathbf{L}_2$ containing a solenoid, are made to interfere on a screen. The magnetism on $\mathbf{L}_2$ is related to the circulation around the boundary $\mathbf{B}_2 = \partial\mathbf{L}_2$ by Stokes's theorem $\mathbf{St}\,\Theta_2$, the second equality of

(10) $$\mathbf{F}_2 := \oint_{\partial\mathbf{L}_2} \mathbf{P}_2 = \iint_{\mathbf{L}_2} d\mathbf{P}_2 = \iint_{\mathbf{L}_2} \mathbf{S}_2.$$

---

[25] See Wallace (2014), especially page 4.

[26] Aharonov & Bohm (1959) page 487: "As yet no direct experiments have been carried out which confirm the effect of potentials where there is no field. It would be interesting therefore to test whether such effects actually exist."

[27] Aharonov & Bohm (1959) page 487: "This effect will exist, even though there are no magnetic forces acting in the places where the electron beam passes."

[28] *Cf.* footnote 11 above.

One might as well define as many loops or concentric circles as are needed to characterise and separate the entities involved as clearly as possible—whatever may figure in the actual experiment, and however zealously one may want to brandish Ockham's razor. We can confine the magnetic field

$$\mathbf{S}_2 = d\mathbf{P}_2 \in \bigwedge\nolimits^2 \mathbb{R}^2$$

produced by the solenoid to an inner disc[29] $\mathbf{I}_2$ bounded by a circle $\partial\mathbf{I}_2$ which can be thought of as a shield 'keeping all the $\mathbf{S}_2$ in.'[30] The only things that don't vanish on $\mathbf{I}_2$ are $\mathbf{S}_2 = d\mathbf{P}_2$ and its primitive, the magnetic potential

$$\mathbf{P}_2 := \mathsf{A} \in \bigwedge\nolimits^1 \mathbb{R}^2.$$

A larger concentric disc $\mathbf{J}_2$ is worth having too, to allow the definition of an annulus $\mathbf{D}_2 := \mathbf{J}_2 - \mathbf{I}_2$ whose *Gedanken* function is to separate $\mathbf{S}_2$ and $\Psi$ as clearly as possible, introducing a nonvanishing radial distance between them. The annulus will be useful in §2.3, where different interpretations get told apart by what they (do or don't) put in it. The circumference $\partial\mathbf{J}_2$ can be seen as a shield 'keeping the wavefunction out.' There's no reason to think of the circles $\partial\mathbf{I}_2$ and $\partial\mathbf{J}_2$ as mere loops or 'topological circles' (like the deformable $\partial\mathbf{L}_2$); since they're not 'integration loops,' they might as well be concentric circles. The circles not only serve to delimit the conceptually useful annulus $\mathbf{D}_2$ occupied only by $\mathbf{P}_2$, but also separate the two impenetrability conditions[31] expressed by Earman's two limits $n = \infty$, $L = \infty$ and brought together in the ideal Hamiltonians $H_{AB}^I$ (Shech) and $\overline{H}_{AB}^{\mathbf{A}_\infty}$ (Earman). Earman's various "more realistic Hamiltonians" $H_{L,n}$ allow magnetic field lines to close, thus forming loops ($L < \infty$); and $\Psi$ to tunnel inwards ($n < \infty$).

Outside the larger disc $\mathbf{J}_2$, the magnetic field $\mathbf{S}_2$ (which already vanished on $\mathbf{D}_2$) still vanishes, but $\mathbf{P}_2$ and the wavefuction don't, thus allowing them to interact. The two different integrals equated by Stokes's theorem require the deformable loop $\partial\mathbf{L}_2$ to be the boundary of a region $\mathbf{L}_2$, which needn't be a disc; but here it does have to be simply-connected (which here just means *unperforated*), and to contain (the support of) $\mathbf{S}_2$, in other words $\mathbf{I}_2 \subset \mathbf{L}_2$. The loop $\partial\mathbf{L}_2$ can otherwise be freely deformed: without affecting the circulation $\mathbf{F}_2$ it can be placed in the annulus $\mathbf{D}_2$, or even outside it; it can even cut the circle $\partial\mathbf{J}_2$ an arbitrary (but even) number of times. But it may be best to leave $\partial\mathbf{L}_2$ outside $\mathbf{J}_2$, where it can be thought of as the 'locus of interaction' between the wavefunction and magnetic potential $\mathbf{P}_2$.

All three loops $\partial\mathbf{L}_2$, $\partial\mathbf{J}_2$, $\partial\mathbf{I}_2$ aren't strictly necessary. Enthusiastic exercise of Ockham's razor could eliminate at least one of them, but intelligibility and conceptual clarity would suffer: the wavefunction and magnetic field would become contiguous, any action at-a-distance would be less evident *etc*.

The $\mathbf{AB}_2$ effect is essentially this: varying the current through the solenoid shifts

---

[29]This corresponds to the $S_{in}$ of Shech (2015, 2018a). Here Shech's $S_{in}$ and $S_{out}$ are separated by the annulus $\mathbf{D}_2$.

[30]In two dimensions it makes sense to speak, somewhat sloppily, of 'keeping the magnetism in'; what one really wants to prevent, however, is magnetism going *above & below* the solenoid, rather than just *through* it. But even an extremely long solenoid would rely on a 'destructive interference' outside $\mathbf{I}_2$ that may not be perfect—a possibility related to the central thesis of Mattingly (2007).

[31]Wallace (2014) page 4: "The conceptual problem is that a sufficiently well-constructed and well-shielded solenoid will result both in negligible magnetic field *outside* the solenoid, and negligible wavefunction *inside* the solenoid."

the interference pattern *differentially*, by the phase factor[32]

$$\exp i \oint_{\mathbf{B}_2} \mathbf{P}_2;$$

small changes in the current barely shift the pattern, larger changes shift it more.[33] This differential character is worth bearing in mind because the topological property on which the *topology interpretation* (§§2.3.4, 3.3.4 below) rests is binary: *perforated or not*. How can a binary property, a mere dichotomy, account for all the nuances of a differential law?

Even if I hesitate to confuse the reader with three-dimensional objects in what I'm sloppily (but rightly) calling the 'two-dimensional' case ($\mathbf{AB}_2$), the magnetic field is in fact a two-form

$$\mathsf{B} = d\mathsf{A} \in \bigwedge{}^{2} \mathbb{R}^3$$

in three dimensions produced by the current density

$$\mathsf{J} = d\!*\!\mathsf{B} = d\mathsf{H} \in \bigwedge{}^{2} \mathbb{R}^3$$

in the solenoid, where the magnetic three-potential

$$\mathsf{A} \in \bigwedge{}^{1} \mathbb{R}^3$$

is a one-form in three dimensions, the Hodge star $*$ turns the two-form $\mathsf{B}$ into the one-form

$$\mathsf{H} = *\mathsf{B} \in \bigwedge{}^{1} \mathbb{R}^3,$$

and the `curl`

$$d : \bigwedge{}^{1} \mathbb{R}^3 \to \bigwedge{}^{2} \mathbb{R}^3$$

turns one-forms into two-forms.

Once we confine our attention to the $xy$ plane, the two-form

$$\mathsf{B} \in \bigwedge{}^{2} \mathbb{R}^2$$

becomes a density (as does $\mathsf{J}$), whereas one-forms in

$$\bigwedge{}^{1} \mathbb{R}^2$$

remain covector fields and the `curl`

$$d : \bigwedge{}^{1} \mathbb{R}^2 \to \bigwedge{}^{2} \mathbb{R}^2$$

now turns one-forms into densities.

This paper isn't really about the various idealisations typically involved in $\mathbf{AB}_2$, which include the following:

---

[32] I just take this standard and very accurate expression for granted. But one can no doubt do even better, see Earman (2017) §5.1 for details.

[33] Aharonov & Bohm (1959) page 487: "Instead, it would be easier to vary the magnetic flux within the same exposure for the detection of the interference patterns. Such a variation would [. . .] alter the sharpness and the general form of the interference bands."

- classical electromagnetism (alongside a quantum-mechanical wavefunction)[34]

- infinitely long solenoid

- impenetrable barrier(s).

Again, idealisations are thoroughly considered in Shech (2015, 2018a) and Earman (2017).

To conclude this description of $\mathbf{AB}_2$ I'll summarise its bare logic in a few words. I take[35] all the quantum mechanics (wavefunction, inteference pattern shift *etc.*) to be *no more than a way of measuring the flux* (10) *through the boundary* $\mathbf{B}_2$ of the region $\mathbf{L}_2$ (in other words the circulation of $\mathbf{P}_2$ around a loop $\mathbf{B}_2$ enclosing the solenoid). $\mathbf{F}_2$ and $\mathbf{S}_2 = d\mathbf{P}_2 = d\bar{\mathbf{P}}_2$ are measurable, as is the class $[\mathbf{P}_2] = d^{-1}\mathbf{S}_2$; but since the individual $\mathbf{P}_2$ isn't, it can be freely deformed by a gauge transformation (5).

## 2.3 Four interpretations in two dimensions

Before turning to $\mathbf{AB}_3$ I'll sketch the four interpretations considered in two dimensions. The first three are distinguished by what they put in the annulus $\mathbf{D}_2$ (outside the inner disc $\mathbf{I}_2$ containing the support of the magnetic field $\mathbf{S}_2$):

$\mathbf{AB}_2$-1: nothing at all

$\mathbf{AB}_2$-2: $[\mathbf{P}_2] = d^{-1}d\mathbf{P}_2$

$\mathbf{AB}_2$-3: homotopy class $\mathbf{H}_2 := [\mathbf{B}_2]$.

### 2.3.1 Magnetic field interpretation $\mathbf{AB}_2$-1

It is hard to deny that $[\mathbf{P}_2] = d^{-1}d\mathbf{P}_2$ is *somehow* present in the annulus $\mathbf{D}_2$—perhaps only *mathematically*. What one can try to deny is that the mathematical presence corresponds to anything genuinely physical. The magnetic field interpretation[36] can be expressed as follows:

How can a whole class $[\mathbf{P}_2]$, from which one wouldn't know how to extract a particular individual—uncontroversially at any rate—correspond to a physical entity capable of propagating an influence across the annulus? Surely there's nothing physically meaningful in the annulus, and the influence somehow gets from the solenoid to the wavefunction or screen *despite* the annulus, without being propagated *through* the annulus by a real medium.[37]

So the annulus gets evacuated, at the stroke of a pen, to keep it from being too full. This is tantamount to denying Buridan's ass all sustenance, preemptively as it were, to prevent starvation from symmetry.

---

[34]The idealisation is unfair: the electrons in the 'interference' beam going around the solenoid are treated quantum-mechanically, the ones running through the solenoid classically; but this is no place to deplore the injustice.

[35]*Cf.* footnote 11.

[36]Aharonov & Bohm (1959) page 490: "we might try to formulate a nonlocal theory in which, for example, the electron could interact with a field that was a finite distance away. Then there would be no trouble in interpreting these results, but, as is well known, there are severe difficulties in the way of doing this." See also Healey (1997).

[37]*Cf.* Mattingly (2006) page 246: "how is it that the information about the state of the field goes from the interior of the solenoid out to the path over which the integral is taken?"

Here the popular `slogan` mentioned in the Introduction means something like: the derivative

$$\mathbf{S}_2 = d(\mathbf{P}_2 + d\lambda_2)$$

is alone real, its primitive $\mathbf{P}_2$ is a physically meaningless fiction. Why? Because alongside a single derivative $\mathbf{S}_2$ there are many primitives $[\mathbf{P}_2 + d\lambda_2]_{\lambda_2}$, and multiplicity is awkward, embarrassing, involves choice *etc.* So 'reality' seems to have to do with the destruction of information, of detail, by differentiation—which helps us make up our minds, or rather spares us the trouble of having to do so. But a given object can be both derivative and primitive: velocity is at once the primitive of acceleration and the derivative of position. Do we then have *degrees* of reality? Velocity is more real than position, less real than acceleration? The more information is destroyed, the more reality is extracted? Is ultimate reality reached once there's nothing left to destroy? One could even propose alternatives to the `slogan`: *kernels cleanse* or perhaps *differentiation consolidates reality*.

Consider the de Rham sequence

$$\zeta_1 \mapsto \zeta_2 = d(\zeta_1 + d\zeta_0)$$
$$\zeta_2 \mapsto \zeta_3 = d(\zeta_2 + d(\zeta_1 + d\zeta_0))$$
$$\zeta_3 \mapsto \zeta_4 = d(\zeta_3 + d(\zeta_2 + d(\zeta_1 + d\zeta_0))),$$

where $\zeta_k$ is a $k$-form. The two-form $\zeta_2$ is *real* at first, being invariant under the deformation $d\zeta_0$, only to become an *unreal* deformation itself. Since forms of different degrees can appear side-by-side in physics, a given object can therefore be both invariant/real and fictitious.[38] If it is objected that an entity has to be *measurable or not*—one has to make up one's mind—that's precisely the 'binary' oversimplification, the black & white measurability that clearly deserves complication by shades of grey ($\alpha_3$). Sometimes the various derivatives are all measurable; take position, velocity, acceleration, *etc.*, in other words the $N$-th derivatives of position, with $N = 0, 1, 2, \ldots$; are they all real? Or is each derivative more real than its primitives, in a *crescendo* of ontological intensity and excitement?[39]

### 2.3.2 Potential interpretation $\mathrm{AB}_2$-2

> [...] the potentials must, in certain cases, be considered as physically effective, even when there are no fields acting on the charged particles. [...] we may retain the present local theory and [...] try to give a further new interpretation to the potentials. In other words, we are led to regard $A_\mu(x)$ as a physical variable.[40]

An equivalence class like $[\mathbf{P}_2] = d^{-1}d\mathbf{P}_2$ can, in the potential interpretation, be physically meaningful and even propagate an influence. Modern physics is full of equivalence classes. Purging physics of all equivalence classes is an overambitious, unreasonable and unrealistic agenda. A wavefunction, for instance, is in fact an equivalence class $[\Phi]$ of functions that differ only on a set of vanishing measure; is the cardinality

---

[38] Caccese (2024)

[39] For an excellent summary of surplus structure and other relevant philosophical issues see the first six pages of Ryckman (2003).

[40] Aharonov & Bohm (1959) pages 490-1; see also Feynman & *al.* (1964) §15-5, Peshkin & Tonomura (1989) page 137: "The AB effect is regarded as experimental evidence for the physical reality of gauge fields."

of $[\Phi]$ enough to rule out the physical relevance of wavefunctions and perhaps even the entities (electrons, planets, the universe *etc.*) they describe, consigning them all to a shady realm of mathematical fictions? Or take barycentric colour theory, which goes back to Newton's *Opticks*: each colour sensation corresponds to infinitely many (convex) combinations of primary colours. Does the very abundance of the combinations condemn them to dismissal? Is the invariant sensation all we have, the light that produced it all garbage?

Or in thermodynamics, each temperature $T$ corresponds to a huge class of mechanical micro-states; should the mechanical states then be ignored, rejected? Is *embarras de richesses* so embarrassing that the riches should all be altogether renounced? Or consider, more specifically, the propagation of heat from the centre of an initially cold ball to its spherical surface. All we know is that the micro-mechanical kinetic energy *somehow* gets from the heated centre to the initially cold spherical surface, thereby heating it; but we know nothing of the micro-mechanical details of the propagation, in the sense that there are infinitely many equivalent possibilities of propagation— reminiscent of the infinitely many possibilities of propagation in the equivalence class $[\mathbf{P}_2]$. Does our unsurmountable ignorance of micro-mechanical details oblige us to renounce that way of understanding heat and its propagation?

So the magnetic 'turbulence' produced in the solenoid is somehow radiated through the annulus as the class $[\mathbf{P}_2]$, unhindered by its cardinality or size or numerosity or superabundance or redundancy. And besides, whyever should one be able to capture this interestingly ambiguous 'turbulence' as unambiguously as desired?—why should the desire, however reasonable or understandable, admit of complete satisfaction? Maybe the interesting ambiguity should just be accepted with resignation and patience.

### 2.3.3 Holonomy interpretation $\mathbf{AB}_2$-3

The equivalence class $[\mathbf{P}_2]$ is somehow more disturbing or more obviously 'infinite' than the homotopy class $\mathbf{H}_2 = [\mathbf{B}_2]$ of loops going around the solenoid once, which is therefore preferred—in the holonomy interpretation—as a medium of propagation.[41] So the '$\mathbf{P}_2$-turbulence' produced at the source $\mathbf{S}_2$, which could manifest itself as '$\mathbf{P}_2$-circulation' at any (arbitrarily distant) loop $\mathbf{B}_2 \in \mathbf{H}_2$, is somehow conveyed by $\mathbf{H}_2$. See also §2.3.5 below.

Again, the first three interpretations differ in their attitudes to the embarrassing riches of the equivalence class $[\mathbf{P}_2]$, which can be found

$\mathbf{AB}_2$-1: so embarrassing that they should be altogether renounced

$\mathbf{AB}_2$-2: not overly embarrassing

$\mathbf{AB}_2$-3: so embarrassing they should be replaced by equally embarrassing riches dual to them—see §2.3.5 below.

---

[41]See Wu & Yang (1975), Healey (1997, 2001, 2004, 2007), Belot (1998), Lyre (2001, 2002, 2004a,b), Myrvold (2011).

### 2.3.4 Topology interpretation $AB_2$-4

The topology interpretation[42] is somewhat different, and seems to rest on a logical mistake[43] along the lines of 'affirming the consequent.'

Here the concentric discs $\mathbf{I} = \mathbf{I}_2$ and $\mathbf{J} = \mathbf{J}_2$ are enough, and it is enough for $\mathbf{I}$ to *contain*[44] the support of the source $\mathbf{S} = \mathbf{S}_2$; we can do without the deformable region $\mathbf{L} = \mathbf{L}_2$. Even if $\mathbf{I}$ and $\mathbf{J}$ are in fact just as deformable as $\mathbf{L}$, it is perhaps easiest to think of them as being rigid for the time being.

Being a disc, $\mathbf{J}$ is simply connected, its only boundary is the circumference $\partial\mathbf{J}$. We can therefore write[45]

$$\mathbf{F}_{\partial\mathbf{J}} = \oint_{\partial\mathbf{J}} \mathsf{A} = \iint_{\mathbf{J}} d\mathsf{A} = \iint_{\mathbf{J}} d(\mathsf{A} + d\lambda_2),$$

where the arbitrary function

$$\lambda_2 \in \bigwedge{}^0 \mathbf{J}$$

is worth bearing in mind (despite having no rôle in this argument), and

$$\mathsf{A} \in \bigwedge{}^1 \mathbf{J}$$
$$d\mathsf{A} = \mathsf{B} \in \bigwedge{}^2 \mathbf{J}.$$

We also have the four implications

(11) $$(\mathsf{B} = 0|_{\mathbf{J}}) \Leftrightarrow (\mathsf{A} \text{ is exact}|_{\mathbf{J}}) \Leftrightarrow (\mathbf{F}_{\partial\mathbf{J}} = 0),$$

where it is assumed, with no (magnetically) interesting loss of generality, that $\mathsf{B} \geq 0$.

By perforating $\mathbf{J}$ with the hole $\mathbf{I} \subset \mathbf{J}$, we obtain a new boundary $\partial\mathbf{I} = -\partial_i\mathbf{D}$ which is external with respect to $\mathbf{I}$ but internal with respect to the annulus $\mathbf{D} = \mathbf{J} - \mathbf{I}$, where the opposite signs, and the directions of rotation they represent, reflect the different perspectives. The 'total' boundary

$$\partial\mathbf{D} = \partial_e\mathbf{D} \cup \partial_i\mathbf{D}$$

---

[42] Shech (2015) page 1078: "the *Aharonov-Bohm* (AB) *effect* [. . . ] is also standardly but falsely dubbed as "topological" and attributed to the presence of a non-simply connected topological space"; Shech (2018b) §5.2: "it is standardly claimed that the AB effect can only occur in a non-simply connected configuration space." See Afriat (2013), especially footnote 6, especially Batterman (2003, pages 544, 552-3, 554-5). *Cf.* Mattingly (2006) page 255: "So, given that $\mathbf{A}$ is non-zero, the integral around the excised region will be non-zero, and there will be a shift in the interference fringes"; even if $\mathbf{A}$ is closed ($\nabla \times \mathbf{A} = 0$) and there's a hole, the integral can vanish.

[43] Shech (2015) page 1081 seems to be making a related point: "it is not legitimate to model the configuration space that corresponds to $H_{AB}^I$ [. . .] as an idealized non-simply connected space. This is so because [. . .] such topological idealizations are pathological in the sense that the non-simply connected topology is a property of a limit system that does not match the corresponding limit property—any minute de-idealization renders the topology of the space simply connected. This is also the case for systems manifesting the AB effect: as long as the solenoid is finite and penetrable, the space in which the AB effect takes place will be simply connected. It is only at the limit that the property of non-simple connectedness arises, so it makes no sense to talk about the space being "approximately" multiply connected."

[44] The two cases (separarated by "whereas") in footnote 10 of Earman (2017) amount to a distinction between loops that go around the solenoid and those that don't; but the distinction is only meaningful if the solenoid is on, otherwise there's just a single homotopy class. Only the loops going around the solenoid are sensitive to the current in the solenoid (and to its magnetic field); but that's just physics, not a topological theorem. The fact that Stokes's theorem cannot be applied within the non-simply connected electron configuration space $\mathcal{R} = \mathbb{R}^3 \backslash \mathcal{S}_\infty$ means that the line integral *may not vanish*, not that *the solenoid is on*.

[45] The subscript $s$ below a flux $\mathbf{F}_s$ is usually a small integer to do with dimensionality, here it will represent the boundary through which the flux is taken; but no confusion is possible.

of the annulus therefore has two parts, an external circle $\partial_e \mathbf{D} = \partial \mathbf{J}$ and an internal circle $\partial_i \mathbf{D} = -\partial \mathbf{I}$. We can still write

$$\mathbf{F}_{\partial \mathbf{D}} = \oint_{\partial \mathbf{D}} \mathsf{A} = \iint_{\mathbf{D}} d\mathsf{A} = \iint_{\mathbf{D}} d(\mathsf{A} + d\lambda_2),$$

but now

$$\oint_{\partial \mathbf{D}} \mathsf{A} = \oint_{\partial_e \mathbf{D}} \mathsf{A} - \oint_{\partial_i \mathbf{D}} \mathsf{A},$$

where the negative sign on the right reflects the different orientations. Of the four implications (11), three survive. The two implications

$$(\mathsf{B} = 0|_{\mathbf{D}}) \Leftrightarrow (\mathbf{F}_{\partial \mathbf{D}} = 0)$$

still hold; and while

$$(\mathsf{A} \text{ is exact}|_{\mathbf{D}}) \Rightarrow (\mathsf{B} = 0|_{\mathbf{D}})$$

remains valid, we now have that

$$(\mathsf{B} = 0|_{\mathbf{D}}) \nRightarrow (\mathsf{A} \text{ is exact}|_{\mathbf{D}}).$$

Assuming that $\mathbf{B}$ vanishes on $\mathbf{D}$ we obtain the equivalence

$$(\mathsf{B} = 0|_{\mathbf{I}}) \Leftrightarrow (\mathsf{A} \text{ is exact}|_{\mathbf{D}}).$$

While we have that

$$(\mathsf{B} = 0|_{\mathbf{J}}) \Leftrightarrow (\mathbf{F}_{\partial \mathbf{J}} = 0)$$
$$(\mathsf{B} = 0|_{\mathbf{D}}) \Leftrightarrow (\mathbf{F}_{\partial \mathbf{D}} = 0) \Leftrightarrow (\mathbf{F}_{\partial \mathbf{J}} = -\mathbf{F}_{\partial \mathbf{I}}),$$

the new topology given by $\mathbf{I}$ (and hence by $\mathbf{D}$) provides, as long as $\mathsf{B}$ vanishes on $\mathbf{D}$, the (obviously equivalent) implications[46]

$$(\mathsf{B} = 0|_{\mathbf{I}}) \Leftrightarrow (\mathbf{F}_{\partial \mathbf{J}} = 0)$$
$$(\mathsf{B} \neq 0|_{\mathbf{I}}) \Leftrightarrow (\mathbf{F}_{\partial \mathbf{J}} \neq 0).$$

While the implication

$$[(\mathsf{B} = 0|_{\mathbf{D}}) \wedge (\mathsf{B} \neq 0|_{\mathbf{I}})] \Rightarrow (\mathbf{F}_{\partial \mathbf{J}} \neq 0)$$

holds, and we can even write

(12) $$(\mathsf{B} = 0|_{\mathbf{D}}) \nRightarrow (\mathbf{F}_{\partial \mathbf{J}} = 0),$$

the 'topological' interpretation seems to rest on the groundless implication

$$(\mathsf{B} = 0|_{\mathbf{D}}) \Rightarrow (\mathbf{F}_{\partial \mathbf{J}} \neq 0),$$

which is by no means equivalent to (12).

Shech (2015, page 1083) seems to be making about the same point in slightly different terms: "the non-trivial holonomy [my $\mathbf{F}_{\partial \mathbf{J}} \neq 0$] arises from the non-flatness

---

[46]Earman (2017) §6: "there are many possible causes hidden within the interior region of the solenoid (including of course a magnetic flux) that could explain the behaviour of the electrons she observes in the exterior."

$[d\mathsf{A} \neq 0]$ of the derivative operator [essentially $\mathsf{A}$] on the principal bundle [with base $\mathbf{J}$ or $\mathbf{D}$], which need not be non-simply connected [$\mathbf{D}$] for nontrivial holonomies"—indeed there's no need to create a hole $\mathbf{I}$ in which to lodge the source, it is enough to put the source in the simply-connected $\mathbf{J}$ itself, so that

$$\neg(d\mathsf{A} = 0|_{\mathbf{J}}).$$

Shech's "non-trivial holonomy" arises from the 'turbulence' $\mathsf{A}$ (produced at the source where $d\mathsf{A} \neq 0$), which eventually reaches the enclosing loop $\partial\mathbf{J}$. But the (partly topological) condition

$$d\mathsf{A} = 0|_{\mathbf{J}}$$

rules out the existence of an '$\mathsf{A}$-producing' source throughout $\mathbf{J}$; the weaker condition

$$d\mathsf{A} = 0|_{\mathbf{D}}$$

is needed to 'make room' for a source (in $\mathbf{I}$), whose 'turbulence' $\mathsf{A}$ would first cross the inner boundary $\partial\mathbf{I}$ (with one sign $\pm$), then the outer boundary $\partial\mathbf{J}$ (with the other sign $\mp$). So the 'derivative operator' $\mathsf{A}$ has to be curved *somewhere* in $\mathbf{J}$ to prevent the holonomy $\mathbf{F}_{\partial\mathbf{J}}$ from vanishing; if it were flat *throughout* $\mathbf{J}$, the holonomy $\mathbf{F}_{\partial\mathbf{J}}$ would vanish:

$$(d\mathsf{A} = 0|_{\mathbf{J}}) \Rightarrow (\mathbf{F}_{\partial\mathbf{J}} = 0)$$
$$(\mathbf{F}_{\partial\mathbf{J}} \neq 0) \Rightarrow \neg(d\mathsf{A} = 0|_{\mathbf{J}}).$$

Summing up, the topology interpretation appears to rest on the mistake

"If $\mathbf{J}$ *may or may not* contain a source there *will be* a flux through $\partial\mathbf{J}$."

which can be correctly rendered as

"A source in $\mathbf{J}$ will produce a flux through $\partial\mathbf{J}$."

or possibly as

"$\mathbf{J}$ may contain a source which—if present—would produce a flux through $\partial\mathbf{J}$."

Shech (2015, 2018a) and Earman (2017) have the right approach to topology ($\mathbf{J}$ *vs.* $\mathbf{D}$), or rather the hole $\mathbf{I}$: they consider what's inside it and what the wavefunction does at its boundary—the hole on its own, apart from its contents, being relatively uninteresting.

Ampère's law $\mathsf{J} = d\mathsf{H}$ provides an interesting `analogy`: Consider two coaxial cylinders in $\mathbb{R}^3$, and the two circles they cut in a generic perpendicular section $\mathbb{R}^2$. If the magnetic one-form

$$\mathsf{H} = *\mathsf{B} \in \bigwedge\nolimits^{1} \mathbb{R}^2$$

is closed everywhere inside the larger circle, we know there's no current

$$\mathsf{J} = d\mathsf{H} \in \bigwedge\nolimits^{2} \mathbb{R}^3$$

through the outer cylinder, and that the circulation of $\mathsf{H}$ around the outer circle vanishes. But if we only know that $\mathsf{H}$ is closed on the annulus between the circles, there may or may not be a current through the inner cylinder, and the circulation of $\mathsf{H}$ around the outer circle may or may not vanish. The magnetism, the circulation of $\mathsf{H}$ around the outer circle, is produced *by the source* (the current in this case), however, not by logically sloppy topological circumlocution.

We even have the characteristic 'differential' freedom

$$[\mathsf{H}] = d^{-1}d(\mathsf{H} + d\tau):$$

much like the magnetic field $\mathsf{B}$ in $\mathbf{AB}_2$, the current is invariant under the addition of a gradient $d\tau$; indeed this Ampèrian `analogy` almost resembles $\mathbf{AB}_2$ *too* much to be of any use. Analogy rests on *difference*, and collapses at *identity*.

Such, then, are the four interpretations in two dimensions; but before turning to three dimensions, let us return to $\mathbf{AB}_2$-2 and $\mathbf{AB}_2$-3, and their relationship.

### 2.3.5  Duality of $\mathrm{AB}_2$-2 and $\mathrm{AB}_2$-3

Since loop deformations in $\mathbf{H}_2$ can be seen as *dual* to gauge transformations (5), there may be a sense in which the holonomy interpretation is 'dual' to the potential interpretation $\mathbf{AB}_2$-2.

The duality between vectors $\mathbf{v}$ in a vector space $\mathbb{V}$ and covectors $\upsilon$ in its dual $\mathbb{V}^*$ is intimately related to the invariance of the scalar product

$$\langle \,\cdot\,,\cdot\, \rangle : \mathbb{V}^* \times \mathbb{V} \to \mathbb{R}$$
$$(\upsilon, \mathbf{v}) \mapsto \langle \upsilon, \mathbf{v} \rangle;$$

one can accordingly consider an even more invariant `scalar`$_2$ product

$$\langle \,\cdot\,,\cdot\, \rangle_2 : \mathbf{H}_2 \times \bigwedge\nolimits^1 \mathbb{R}^2 \to \mathbb{R}$$
$$(\mathbf{B}_2, \lambda_2) \mapsto \langle \mathbf{B}_2, \lambda_2 \rangle_2 := \oint_{\mathbf{B}_2} (\mathbf{P}_2 + d\lambda_2),$$

with the important difference that covectors in $\mathbb{V}^*$ are *naturally* paired $\mathbb{V}^* \leftrightarrow \mathbb{V}$ with vectors in $\mathbb{V}$. In general relativity, for instance, the metric $g_{\mu\nu}$ achieves the pairing between the vector components and the dual components

$$V^\nu = \sum_\mu g^{\mu\nu} V_\mu$$

very naturally by 'lowering the index,' thus producing an isomorphism between the tangent and cotangent spaces. Even if I can't think of a *natural* way of pairing the deformations

- of loops $\mathbf{B}_2 \in \mathbf{H}_2$

- $(\lambda_2)$ of potentials/primitives $\mathbf{P}_2 \in [\mathbf{P}_2]$,

there are no doubt ways of defining rather unnatural pairings

$$\lambda_2 : \mathbf{H}_2 \to \mathbf{H}_2$$
$$\mathbf{B}_2 \mapsto \bar{\mathbf{B}}_2(\lambda_2)$$

by giving the deformed loop $\bar{\mathbf{B}}_2(\lambda_2)$ this or that dependence on the function $\lambda_2$ serving to deform the potential. The most 'natural' pairing between $\mathbf{H}_2$ and $[\mathbf{P}_2]$ that comes to mind may be a single diffeomorphism that serves two rather different purposes: on the one hand it deforms loops $\mathbf{B}_2 \in \mathbf{H}_2$ by dragging them, point by point; and on

the other it (or rather its tangent map) drags the potential $\mathbf{P}_2 \in [\mathbf{P}_2]$. The diffeomorphism achieves the pairing by assigning a loop deformation to every deformation of the potential.

The invariance of the scalar product $\langle\,\cdot\,,\cdot\,\rangle$ depends entirely on natural pairing, which naturally assigns a covector $\upsilon_{\mathbf{v}} =: \mathbf{v}^\flat \in \mathbb{V}^*$ to every vector $\mathbf{v} \in \mathbb{V}$; dual representations of the same group, say $\mathbf{SO}(N)$, guarantee the invariance by making one basis (and its components) co-vary while the dual basis contra-varies: the two variations are co-ordinated in just the right way, for perfect compensation. But the $\text{scalar}_2$ product $\langle\,\cdot\,,\cdot\,\rangle_2$ will remain invariant however $\mathbf{B}_2 \in \mathbf{H}_2$ gets deformed, *and* however $\mathbf{P}_2 \in [\mathbf{P}_2]$ gets deformed by $\lambda_2$; even if the invariance would not be *upset* by a pairing, it would by no means *rely* on it in the way the invariance of the scalar product $\langle\,\cdot\,,\cdot\,\rangle$ relies on the natural pairing $\mathbf{v} \mapsto \upsilon_{\mathbf{v}} = \mathbf{v}^\flat$.

Even if the duality between $\mathbf{H}_2$ and $[\mathbf{P}_2]$ is accordingly 'weaker' than other dualities, since the corresponding invariance doesn't rely on a natural pairing, we can still view the holonomy interpretation $\mathbf{AB}$-3 as 'dual'—perhaps even *equivalent*—to the potential interpretation $\mathbf{AB}$-2; suggesting not only that the potential interpretation may be *no worse than* the holonomy interpretation, but even the question: Why go to all the trouble of formulating the dual interpretation, if nothing is thereby gained? I won't press the point, because the duality itself may be interesting enough to deserve formulation—but only as an alternative worth bearing in mind, not as a solution to the problem of choice posed by an equivalence class: it only replaces one embarrassment with another.

For a rough understanding of the duality, the 'lines' of $[\mathbf{P}_2]$ can be crudely imagined as 'across' the turbulence, those of $\mathbf{H}_2$ as 'along' it; a representative $\mathbf{P}_2 \in [\mathbf{P}_2]$ (one can think, for instance, of straight lines radiating out from the solenoid) would accordingly be perpendicular to a representative $\mathbf{B}_2 \in \mathbf{H}_2$ (one can think of a circle centred on the solenoid).

Needless to say, the same duality holds in three dimensions.

# 3  Three dimensions

## 3.1  $\mathbf{AB}_3$

In three dimensions we have a mass/charge density

$$\text{(13)} \qquad\qquad \mathbf{S}_3 := \rho = dE \in {\bigwedge}^3 \mathbb{R}^3$$

radiating a gravitational/electrostatic field

$$\mathbf{P}_3 := E \in {\bigwedge}^2 \mathbb{R}^3$$

which is then caught by an enclosing membrane $\mathbf{B}_3 = \partial\mathbf{L}_3$; $\mathbf{St}\,\Theta_3$ is the second of the three equalities

$$\mathbf{F}_3 = \iint_{\partial\mathbf{L}_3} \mathbf{P}_3 = \iiint_{\mathbf{L}_3} d\mathbf{P}_3 = \iiint_{\mathbf{L}_3} \mathbf{S}_3,$$

the third being the integral version of the differential Gauss-Maxwell electrostatic law (13). The 'three-dimensional annulus' or 'cavity' $\mathbf{D}_3$ is the difference $\mathbf{D}_3 = \mathbf{J}_3 - \mathbf{I}_3$ between the (spherical) support $\mathbf{I}_3$ of $\mathbf{S}_3$ and the larger concentric ball $\mathbf{J}_3$, with spherical boundary $\partial\mathbf{J}_3$. For the analogy with $\mathbf{AB}_2$ to hold we need a three-dimensional

freedom along the lines of

$$[\mathbf{P}_2] = [\mathbf{P}_2 + d\lambda_2]_{\lambda_2} = d^{-1}\mathbf{S}_2,$$

where

$$\lambda_2 \in \bigwedge\nolimits^{0} \mathbb{R}^2$$

is a function and $[\mathbf{P}_2]$ the equivalence class of all one-forms differing by an exact term $d\lambda_2$, which, much like a constant of integration, picks a one-form out of the kernel of the `curl`

$$d : \bigwedge\nolimits^{1} \mathbb{R}^2 \to \bigwedge\nolimits^{2} \mathbb{R}^2.$$

In a sense we already have the right freedom in

$$[\mathbf{P}_3] = [\mathbf{P}_3 + d\lambda_3]_{\lambda_3} = d^{-1}\mathbf{S}_3,$$

where

$$d\lambda_3 \in \bigwedge\nolimits^{2} \mathbb{R}^3$$

is the `curl` of a one-form

$$\lambda_3 \in \bigwedge\nolimits^{1} \mathbb{R}^3.$$

Again, the difference is that the gravitational/electrostatic field $\mathbf{P}_3$ is usually considered more measurable than the magnetic potential $\mathbf{P}_2$; but in §3.2.1 below I'll take the individual field $\mathbf{P}_3$ to be about as unmeasurable as $\mathbf{P}_2$ in $\mathbf{AB}_2$.

The density $\mathbf{S}_3$, flux $\mathbf{F}_3$ and equivalence class $[\mathbf{P}_3]$ are all measurable. Since $\mathbf{F}_3$ and $\mathbf{S}_3$ are related by Stokes's theorem, a gradual increase—however arranged or conceived ($O_2$)—in $\mathbf{S}_3$ would produce a corresponding increase in $\mathbf{F}_3$.

It is now easier to see why $\mathsf{B}$ is what I've called a *source* $\mathbf{S}$. Just as the mass/charge density $\mathbf{S}_3$ is the source of the gravitational/electrostatic field $E = \mathbf{P}_3$ (a perturbation, an alteration of the surrounding medium) which is then caught by the boundary $\mathbf{B}_3$, the magnetic field $\mathsf{B} = \mathbf{S}_2$ is the source of the surrounding magnetic 'turbulence' which, conveyed by the potential $\mathsf{A} = \mathbf{P}_2$, is likewise caught by the boundary $\mathbf{B}_2$.[47]

I'll conclude this description with a few words about the bare logic of $\mathbf{AB}_3$. Since $\mathbf{AB}_2$ involved $\mathbf{St}\,\Theta_2$, $\mathbf{AB}_3$ will involve

$$\mathbf{F}_3 = \iint_{\partial\mathbf{L}_3} \mathbf{P}_3 = \iiint_{\mathbf{L}_3} d\mathbf{P}_3 = \iiint_{\mathbf{L}_3} \mathbf{S}_3.$$

For the analogy to work, the flux $\mathbf{F}$ through the boundary $\mathbf{B} = \partial\mathbf{L}$ and the source $\mathbf{S}$ (here $\mathbf{F}_3$ through $\mathbf{B}_3 = \partial\mathbf{L}_3$, and $\mathbf{S}_3 = d(\mathbf{P}_3 + d\lambda_3)$) have to be just as measurable as in $\mathbf{AB}_2$; and the primitive $\mathbf{P}$ (here the gravitational/electrostatic field $\mathbf{P}_3 := E$) has to be about as unmeasurable as in $\mathbf{AB}_2$. The unmeasurability of $\mathbf{P}_3$ yields the freedom

(14)     $$\mathbf{P}_3 \mapsto \bar{\mathbf{P}}_3 = \mathbf{P}_3 + d\lambda_3$$

and the equivalence class $[\mathbf{P}_3] = d^{-1}d\mathbf{P}_3$.

Poisson's equation, which we've managed to do without so far, will appear in the next section.

---

[47] *Cf.* the `analogy` at the end of §2.3.4 above.

## 3.2 Unmeasurability and measurability ($\alpha_2$)

### 3.2.1 Just how measurable is the electrostatic field?

How does one get from the Gauss-Maxwell law

$$(\mathbf{GM}) \qquad\qquad d\mathbf{P}_3 = d\bar{\mathbf{P}}_3 = d(\mathbf{P}_3 + d\lambda_3) = \mathbf{S}_3,$$

with all its potentially troublesome gauge freedom (14), to Coulomb's spherically symmetrical law **CL**? To get from **GM** to Poisson's equation $d*d\phi = \mathbf{S}_3$ one needs the exactness condition

$$(\mathbf{EX}_3) \qquad\qquad \mathbf{P}_3 = *d\phi,$$

which is locally equivalent to the condition $d*\mathbf{P}_3 = 0$ derived, provided $\partial_t\mathsf{B}$ vanishes (as it does in magnetostatics), from the Maxwell-Faraday equation[48]

$$d*\mathbf{P}_3 = -\partial_t\mathsf{B}.$$

And to get from Poisson's equation to **CL**, asymptotic constancy[49] **AC** is needed:

$$\mathbf{GM}\,\&\,\mathbf{EX}_3\,\&\,\mathbf{AC} \Rightarrow \mathbf{CL}.$$

So it takes quite a lot—$\mathbf{EX}_3\,\&\,\mathbf{AC}$—to fix gauge (**CL**) in $\mathbf{AB}_3$. The two assumptions are in fact independent: the asymptotic assumption can also be referred to the force ($\mathbf{P}_3 = *d\phi \to 0$) rather than its potential. Both assumptions are needed, since either one on its own would leave much freedom. **AC** alone could be seen as applying beyond a sphere of finite radius, inside which there would remain the freedom

$$\mathbf{P}_3 \mapsto \bar{\mathbf{P}}_3 := \mathbf{P}_3 + d\xi_3$$

to deform with an exact term.

Or the principle **PSR** of sufficient reason can be invoked to provide the rotational symmetry (equivariance) of Coulomb's law **CL**, thus overcoming all the gauge freedom left by the Gauss-Maxwell law **GM**. But the **PSR** would require a perfect homogeneity that never holds exactly, being undermined by (generally nontrivial) constitutive relations and/or curvature. The dielectric constant $\varepsilon$ sometimes figures explicitly in Coulomb's law—but indeed *as a constant*. If we allow $\varepsilon$ to vary, Coulomb's law ceases to make sense; and makes even less sense if $\varepsilon$ becomes a tensor, to describe an anisotropic medium. And besides, in ideal circumstances the **PSR** could be used to fix gauge in $\mathbf{AB}_2$ as well, picking out a single equivariant element

$$\mathbf{P}_{\mathrm{PSR}} \in [\mathbf{P}] = d^{-1}\mathbf{S}.$$

Or how about test bodies—surely *they* can fix the field in **NG/ES**. If test bodies are *Gedanken* fantasies, in limitless supply, maybe we can help ourselves to them in $\mathbf{AB}_2$

---

[48]This is not the place to justify $\mathbf{EX}_3$ in **NG**.

[49]In **NG**, Einstein (1917, page 142) just takes $\mathbf{EX}_3$ for granted (by writing Poisson's equation) but assumes **AC** explicitly: "Es ist wohlbekannt, daß die POISSONsche Differentialgleichung [. . .] in Verbindung mit der Bewegungsgleichung des materiellen Punktes die NEWTONsche Fernwirkungstheorie noch nicht vollständig ersetzt. Es muß noch die Bedingung [**AC**] hinzutreten, daß im räumlich Unendlichen das Potential $\phi$ einem festen Grenzwerte zustrebt"—finding 'asymptotic' assumptions such as constancy (**AC**) or flatness philosophically unsatisfactory and disturbingly unMachian, however. Interestingly, he never says a word about test bodies.

as well, to measure $\mathbf{P}_2$ somehow or other. If they're concretely experimental instead, and subject to all the unfortunate limitations of the real world, they won't be too numerous; and how can a triply non-denumerable infinity $\infty^3$ of field values $\mathbf{P}_x, \mathbf{P}_y, \mathbf{P}_z$ be uniquely determined by a handful of test bodies (amid nontrivial constitutive relations and/or curvature)? And surely *no* observed values and *a handful* of such values are about equivalent (equinumerous) with respect to the triple non-denumerable infinity;

$$(15) \qquad\qquad (\text{handful} : \infty^3) \approx (0 : \infty^3).$$

It seems significant that Einstein (1917) never says a word about test bodies.

But even if one feels that *a handful of points can constrain the field an awful lot*, and chooses not to take the unmeasurability of $\mathbf{P}_3$ seriously, my scheme is not thereby invalidated, it remains intact—serving as it does to provide a new perspective, shed light *etc*. Besides, one's always free to *assume* that $\mathbf{P}_3$ is unmeasurable, and fruitfully explore the implications; the history of philosophy—even philosophy of physics—being full of enlightening counterfactual assumptions and arguments.

We can conclude this part with the following 'alternative':

- if it makes sense to write (15), $\mathbf{AB}_2$ is *no more teratological than* $\mathbf{AB}_3$

- if $(\text{handful} : \infty^3)$ is to be viewed as *significantly greater than* $(0 : \infty^3)$, the peculiarity of $\mathbf{AB}_2$ would somehow rest on the 'strong inequality'

$$(\text{handful} : \infty^3) \gg (0 : \infty^3),$$

which is surprising enough in itself.

### 3.2.2 Will the magnetic potential remain forever unmeasurable?

In fact there are two ways of 'closing the measurability gap' separating $\mathbf{P}_2$ and $\mathbf{P}_3$, to reinforce the isomorphism $\mathbf{AB}_2 \sim \mathbf{AB}_3$:

- As we've just seen in §3.2.1, one can 'pull $\mathbf{P}_3$ towards $\mathbf{P}_2$' as it were, by showing how problematic the complete, unambiguous individuation of $\mathbf{P}_3$ can be.

- Alternatively one can 'act on $\mathbf{P}_2$' instead—'pushing it towards $\mathbf{P}_3$' as it were, by entertaining the possibility that $\mathbf{P}_2$ may somehow become measurable.

The possible measurability of the (electro)magnetic potential would in itself deserve at least a paper; but I'll do no more than quote Aharonov & Bohm (1959, pages 490-1), who conclude theirs by suggesting that $\mathbf{P}_2$ may not be so unmeasurable after all; quite apart from the actual feasibility of their proposal, the very fact that they even make it is a courageous step in the right direction, of exploring such possibilities, so often hastily dismissed *a priori*.

> [...] we may retain the present local theory and, instead, we may try to give a further new interpretation to the potentials. In other words, we are led to regard $A_\mu(x)$ as a physical variable. This means that we must be able to define the physical difference between two quantum states which differ only by a gauge transformation. It will be shown in a future paper that in a system containing an undefined number of charged particles (i.e., a superposition of states of different total charge), a new Hermitian operator, essentially an angle variable, can be introduced, which is conjugate

to the charge density and which may give a meaning to the gauge. Such states have actually been used in connection with recent theories of superconductivity and superfluidity and we shall show their relation to this problem in more detail.

## 3.3 The four interpretations in three dimensions

The following expressions can be used for the four interpretations in three dimensions:

**AB**$_3$-1: **Gr/Es** *source interpretation*

**AB**$_3$-2: **Gr/Es** *field interpretation*

**AB**$_3$-3: *membrane interpretation*

**AB**$_3$-4: *topology interpretation.*

For the isomorphism **AB**$_3$ $\sim$ **AB**$_2$ captured in the above **Table** to hold it is assumed that $\mathbf{P}_3$ is about as unmeasurable as $\mathbf{P}_2$.

### 3.3.1 Gr/Es source interpretation AB$_3$-1

Since $\mathbf{P}_3$ is assumed unmeasurable, all we have is the class $[\mathbf{P}_3] = d^{-1}d\mathbf{P}_3$, which according to **AB**$_3$-1 has to be a physically meaningless mathematical fiction. So there's nothing physical, nothing real between the source and the boundary $\mathbf{B}_3$. The flux $\mathbf{F}_3$ through the boundary is therefore, in **AB**$_3$-1, a non-local effect, in the sense that it isn't conveyed by a physical 'carrier' $\mathbf{P}_3$.

### 3.3.2 Gr/Es field interpretation AB$_3$-2

In **AB**$_3$-2, the flux $\mathbf{F}_3$ is carried from the source to the boundary $\mathbf{B}_3$ by the gravitational field $\mathbf{P}_3$. To the question "which $\mathbf{P}'_3 \in [\mathbf{P}_3] = d^{-1}d\mathbf{P}_3$ in particular?" there are at least three possible answers:

*a*) It really doesn't matter, any $\mathbf{P}'_3 \in [\mathbf{P}_3]$ will do—all elements of $[\mathbf{P}_3]$ being on an equal footing.

*b*) The elements of $[\mathbf{P}_3]$ are not all on an equal footing; only one of them, $\mathbf{P}'_3$, is the *right one*. But since the distinguished element $\mathbf{P}'_3$ is assumed to be empirically inaccessible, any element of $[\mathbf{P}_3]$ will do.

*c*) The elements of $[\mathbf{P}_3]$ are all on an equal footing, *empirically*. But measurement isn't the only way of selecting or ruling out elements of $[\mathbf{P}_3]$: some could be æsthetically or pragmatically privileged; simplicity, elegance, beauty, economy, convenience, computational considerations or even history[50] could be relevant.

---

[50]Duhem (1989) pages 388ff

### 3.3.3 Membrane interpretation $\mathbf{AB}_3$-3

Since $[\mathbf{P}_3]$ is a class, full of individuals, it has to be—according to $\mathbf{AB}_3$-3—a physically meaningless mathematical fiction. But to avoid the non-locality of $\mathbf{AB}_3$-1, *something* in $\mathbf{D}_3$ has to carry the flux $\mathbf{F}_3$ from the source to the boundary. Since the flux is the same for the whole homotopy class $\mathbf{H}_3$, we can replace $[\mathbf{P}_3]$ with $\mathbf{H}_3$; so a class of boundaries somehow conveys, in $\mathbf{AB}_3$-3, the gravitational/electrostatic field, or rather $[\mathbf{P}_3]$ (or whatever is produced by the source $\mathbf{S}_3$ and then manifests itself at the boundary as a gravitational/electric flux $\mathbf{F}_3$). In other words: There has to be *something* in $\mathbf{D}_3$; if it cannot be the class $[\mathbf{P}_3]$ (being a class), it has to be the class $\mathbf{H}_3$.

Summing up, these first three interpretations can be distinguished by what they put in the isolating region $\mathbf{D}_3$;

$\mathbf{AB}_3$-1 : nothing at all

$\mathbf{AB}_3$-2 : gravitational/electrostatic field

$\mathbf{AB}_3$-3 : homotopy class $\mathbf{H}_2$.

### 3.3.4 Topology interpretation $\mathbf{AB}_3$-4

Most of §2.3.4 applies just as well to $N$ (for instance three) dimensions as to two, so little need be added here.

If there's no source in $\partial\mathbf{J} = \partial\mathbf{J}_N$, the flux $\mathbf{F}$ through $\partial\mathbf{J}$ will vanish; if there may be a source in $\partial\mathbf{J}$, the flux $\mathbf{F}$ through $\partial\mathbf{J}$ may not vanish. The subscript 3 of course gives the three-dimensional case, with $\partial\mathbf{J}_3$ *etc.*

But the flux is undeniably produced by the source, not by the hole itself.[51] Without a source, the implication

$$(d\mathbf{P} = 0|_{\mathbf{D}}) \Rightarrow (\mathbf{F} \neq 0)$$

on which the topological interpretation seems to rest is groundless, indeed wrong.[52] If one cannot avoid attributing the flux to the source that produced it, why bother with all the confusing topological circumlocution? The 'topologist' seems to want to shift the explanation from the source to the topology, replacing the source with an appropriate topological condition. There's nothing wrong with an emphasis[53] on the flux over the source; but a mere hole on its own doesn't even provide the flux. What's to be gained by attributing a planet's gravitational field not to the planet itself but to a possibly empty hole?

## 4   Final remarks

Aharonov & Bohm (1959) view the quantum interaction law (9) as indicating the *fundamentally* (as opposed to *ostensibly*) canonical nature of quantum theory. Even if one can choose to take issue with them on how 'Hamiltonian' quantum mechanics really

---

[51]Earman (2017) §5.3: "the non-simple connectedness plays no causal role in producing the observed effects."

[52]See Afriat (2013).

[53]*Cf.* Weyl (1931) pages 54-5: "Auch bin ich überzeugt, daß die Masse von Hause aus weder träge noch schwere, sondern gravitationsfelderzeugende Masse ist und darum als der Fluß definiert werden muß, den das Gravitationsfeld durch eine durch das Teilchen umschließende Hülle hindurchschickt, so wie nach FARADAY die Ladung der elektrische Kraftfluß durch eine solche Hülle ist."

is, (9)—or more precisely the elimination of (7)—is indeed surprising and interesting enough to deserve the experimental verification they propose. But that's just background, and not what's really at issue in this paper. We've since lost sight of their original motivations, and treated the experiment they propose to test (9) as evidence of *other* peculiarities, to do with non-locality and/or loops and/or topology—which is what this paper is really about. But whatever the rôle and importance of those peculiarities, one certainly can't claim that the Aharonov-Bohm effect has *absolutely nothing* to do with deformable boundaries or nonlocality or topology, which would be going a bit far. So what can one claim? That the effect is, say, *somewhat* topological, but *not too* topological? Hard to formulate and calibrate, to say the least.

So what I've argued is that $\mathbf{AB}_2$ is *no more teratological* (topological or holonomic or whatever) than $\mathbf{AB}_3$; than plain, ordinary, down-to-earth theories that don't seem teratological at all: electrostatics, Newtonian gravity. The strategy resembles a relative consistency proof, inasmuch as I show one thing to be *no worse* than another.

The basic and very general mathematical fact in the background is the (nontrivial) kernel of differential operators, which gives rise to equivalence classes such as $[\mathbf{P}] = d^{-1}d\mathbf{P}$ ($\alpha_2$). The issue ($\alpha_3$) is then one of elimination, selection, even individuation: is one stuck with the embarrassing riches of the whole class $[\mathbf{P}]$, or are there ways of thinning it out, perhaps reducing it to a single element? This is in general an extremely intricate matter, which has to be considered case by case; even in the simplest there can be disagreement as to the available means of elimination, empirical or other. Equivalence classes $[\mathbf{P}]$ are often dismissed as physically meaningless on account of the difficulties involved in reduction to a more manageable size; but such reductions are typically (even in the three-dimensional case $[\mathbf{P}_3]$) problematic. The isomorphism $\mathbf{AB}_2 \sim \mathbf{AB}_3$ I propose is based on a consideration (especially in §3.2) of measurement difficulties, of the 'measurability gap' between $\mathbf{P}_2$ and $\mathbf{P}_3$, which may be smaller than expected.

Measurability is sometimes viewed in crudely 'binary' terms of black & white: this (perhaps the electrostatic field $\mathbf{P}_3$, or even the magnetic field $\mathbf{S}_2 = d\mathbf{P}_2$) *is* measurable whereas that (perhaps the magnetic potential $\mathbf{P}_2$) *isn't*, that's all there is to it. But one can wonder ($\alpha_3$) whether measurability really is as simple as that, whether black & white are enough to capture all the complexities of measurement; 'measurability nuances,' shades of grey, may well be worth interposing between the black and the white. Such shades would then affect our whole understanding of gauge freedom ($\alpha_2$); for instance, just how free are we to deform the gravitational field $\mathbf{P}_3$ with an exact term $d\lambda_3$?

Much turns here on the magnetic gauge freedom $[\mathbf{P}_2] = d^{-1}d\mathbf{P}_2$, where the 'source' $\mathbf{S}_2 = d\mathbf{P}_2$ is supposed to be much more measurable than the potential $\mathbf{P}_2$. In §3.2.1 we saw, however, that $\mathbf{P}_2$'s gravitational/electrostatic analog $\mathbf{P}_3$ is subject to a formally identical freedom $[\mathbf{P}_3] = d^{-1}d\mathbf{P}_3$, which is about as hard to overcome. If the loudly proclaimed weirdness of $\mathbf{AB}_2$ rests on the unmeasurability of $\mathbf{P}_2$, it may therefore be questionable—relying as it does on the insignificant or even dubious difference between the quantities

$$(\text{handful} : \infty^3) \text{ and } (0 : \infty^3)$$

considered in §3.2.1. Aharonov & Bohm (1959) even suggest that the potential $\mathbf{P}_2$ may itself be measurable, which would already be enough to close the 'measurability gap' between $\mathbf{P}_2$ and $\mathbf{P}_3$.

If only one could single out an individual primitive $\mathbf{P} \in [\mathbf{P}]$, *it* would uncontroversially carry the effect from the source to the boundary, without non-locality—or

the need to seek, beyond $\mathbf{P}$, a supposedly 'more invariant' medium (a homotopy class of boundaries) to carry the effect in its place. But even if such straightforward individuation isn't possible, one needn't despair; why should our pervasive measurement difficulties (which afflict $[\mathbf{P}_3]$ as well as $[\mathbf{P}_2]$) and ignorance prevent the class $[\mathbf{P}]$, or perhaps a dimly identifiable part of it, from conveying the effect from the source to the boundary?

As to the topology interpretation, it doesn't even seem to rely on the unmeasurability of the primitive $\mathbf{P}$; so that $\mathbf{AB}_2$ resembles $\mathbf{AB}_3$ even more: $\mathbf{AB}_2$ really is *no more topological than* $\mathbf{AB}_3$, without qualification now. And one wonders how a flux can be due to a mere hole, which may or may not contain a source $\mathbf{S}$; it is clearly produced by the source itself.

The physics of the analogy, if taken too literally, can be misleading: in $\mathbf{AB}_2$, the magnetic field $\mathsf{B}$ is viewed as a 'source' $\mathbf{S}_2 = d\mathbf{P}_2$; in $\mathbf{AB}_3$, the gravitational field $\mathbf{P}_3$ corresponds to the potential $\mathbf{P}_2$ in $\mathbf{AB}_2$; and so on. The analogy makes sense at a more abstract level. Inevitable physical differences between the sources $\mathbf{S}_2$ and $\mathbf{S}_3$, or between the primitives $\mathbf{P}_2$ and $\mathbf{P}_3$, do not obstruct the isomorphism $\mathbf{AB}_2 \sim \mathbf{AB}_3$.

The move to three dimensions produces not one but two new schemes (columns in the above **Table**): an abstract scheme $\mathbf{AB}_N$ in which the essentials of what a logician might call the 'theory' are laid bare, disentangled from the accidental features that would then be confined to the 'models'; and a new three-dimensional model $\mathbf{AB}_3$, to be considered alongside the standard two-dimensional model $\mathbf{AB}_2$. One may even want to view the third column as *subordinate* to the second, in the sense that the best way to extract a theory from a single model can be the identification of a second model to begin with, so one can then see what theory they both satisfy. At any rate the three-dimensional model $\mathbf{AB}_3$ proposed here fruitfully associates *new* accidental features with the essential structures represented in the abstract column $\mathbf{AB}_N$. But as the three-dimensional scheme $\mathbf{AB}_3$ can be considered alongside the other two columns, its accidental features are not *confusingly* entangled with the bare theory $\mathbf{AB}_N$—they can only provide the advantages associated with, for instance, the intuitive familiarity of gravity (a falling apple being more familiar and accessible than a shifting interference pattern) or electrostatics. Such advantages are encountered in metaphors or similes, perhaps even in allegory.

Again, I'm only claiming the three-dimensional analog $\mathbf{AB}_3$ *sheds light* on $\mathbf{AB}_2$, not that it *captures absolutely everything*; indeed capturing everything is exactly what it's meant not to do, that's the whole point—if it did capture every last detail the whole scheme would be sterile and pointless.

# References

Afriat, A. (2013) "Topology, holes and sources" *International Journal of Theoretical Physics* **52**, 1007-12

Aharonov, Y. & D. Bohm (1959) "Significance of electromagnetic potentials in the quantum theory" *Physical Review* **115**, 485-91

Batterman, R. (2003) "Falling cats, parallel parking and polarized light" *Studies in History and Philosophy of Modern Physics* **34**, 527-57

Belot, G. (1998) "Understanding electromagnetism" *The British Journal for the Philosophy of Science* **49**, 531-55

Born, M. (1925) *Vorlesungen über Atommechanik*, Springer, Berlin

Caccese, E. (2024) private communication, *Tre fontane* (Pisciotta), 21 July 2024

Duhem, P. (1989) *La théorie physique: son objet – sa structure*, Vrin, Paris

Earman, J. (2017) "The role of idealizations in the Aharonov-Bohm effect" *Synthese* DOI 10.1007/s11229-017-1522-9

Ehrenberg, W. & R. Siday (1949) "The refractive index in electron optics and the principles of dynamics" *Proceedings of the Physical Society B* **62**, 8-21

Einstein, A. (1917) "Kosmologische Betrachtungen zur allgemeinen Relativitätstheorie" *Sitzungsberichte der Königlich Preussischen Akademie der Wissenschaften*, 142-52

Feynman, R. & *al*. (1964) *The Feynman lectures on physics, volume 2*, Addison-Wesley, Reading

Franz, W. (1965) "Elektroneninterferenzen im Magnetfeld" *Zeitschrift für Physik* **184**, 85-91

Franz, W. (1940) "Elektroneninterferenzen im Magnetfeld" *Physikalische Berichte* **21**, 686

Franz, W. (1939) "Elektroneninterferenzen im Magnetfeld" *Verhandlungen der Deutschen Physikalischen Gesellschaft* **20**, 65-6

Healey, R. (2007) *Gauging what's real: the conceptual foundations of contemporary gauge theories*, Oxford University Press, New York

Healey, R. (2004) "Gauge theories and holisms" *Studies in History and Philosophy of Modern Physics* **35**, 619-42

Healey, R. (1997) "Nonlocality and the Aharonov-Bohm effect" *Philosophy of Science* **64**, 18-41

Healey, R. (2001) "On the reality of gauge potentials" *Philosophy of Science* **68**, 432-55

Hiley, B. (2013) "The early history of the Aharonov-Bohm effect" arXiv:1304.4736v1 [physics.hist-ph]

Kottler, F. (1922) "Maxwell'sche Gleichungen und Metrik" *Sitzungsberichte d. mathem.-naturw. Kl., Abt. IIa (Sitzung 23. Februar)* **131**, 119-46

Lyre, H. (2004a) *Lokale Symmetrien und Wirklichkeit: eine Naturphilosophische Studie über Eichtheorien und Strukturenrealismus*, Mentis, Paderborn

Lyre, H. (2004b) "Holism and structuralism in $U(1)$ gauge theory" *Studies in History and Philosophy of Modern Physics* **35**, 643-70

Lyre, H. (2002) "Zur Wissenschaftstheorie moderner Eichfeldtheorien" in A. Beckermann and C. Nimtz (editors) *Argument & Analyse – Sektionsvorträge*, Mentis, Paderborn

Lyre, H. (2001) "The principles of gauging" *Philosophy of Science* **68**, S371-81

Mattingly, J. (2007) "Classical fields and quantum time-evolution in the Aharonov-Bohm effect" *Studies in History and Philosophy of Modern Physics* **38**, 888-905

Mattingly, J. (2006) "Which gauge matters?" *Studies in History and Philosophy of Modern Physics* **37**, 243-62

Myrvold, W. (2011) "Nonseparability, classical and quantum" *The British Journal for the Philosophy of Science* **62**, 417-32

Olariu, S. & I. Popescu (1985) "The quantum effects of electromagnetic fluxes" *Reviews of Modern Physics* **57**, 339-436

Peshkin, M. & A. Tonomura (1989) *The Aharonov-Bohm effect*, Springer, Berlin

Poincaré, H. (1917) *La science et l'hypothèse*, Flammarion, Paris

Ryckman, T. (2003) "Surplus structure from the standpoint of transcendental idealism: the "world geometries" of Weyl and Eddington" *Perspectives on Science* **11**, 76-106

Shech, E. (2018a) "Idealizations, essential self-adjointness and minimal model explanation in the Aharonov-Bohm effect" *Synthese* **195**, 4839-63

Shech, E. (2018b) "Infinitesimal idealization, easy road nominalism, and fractional quantum statistics" *Synthese*

Shech, E. (2015) "Two approaches to fractional statistics in the quantum Hall effect: idealizations and the curious case of the anyon" *Foundations of Physics* **45**, 1063-1100

Singer, S. (2004) *Symmetry in mechanics: a gentle, modern introduction* Birkhäuser, Boston

Guillemin, V. & S. Shlomo (1984) *Symplectic techniques in physics*, Cambridge University Press

Vaidman, L. (2012) "Role of potentials in the Aharonov-Bohm effect" *Physical Review A* **86**, 040101

Wallace, D. (2014) "Defeating the Aharonov-Bohm effect" arXiv:1407.5073v1

Weyl, H. (1989) *Philosophie der Mathematik und Naturwissenschaft*, Oldenbourg, Munich

Weyl, H. (1931) "Geometrie und Physik" *Die Naturwissenschaften* **19**, 49-58

Wu, T. & C. Yang (1975) "Concept of nonintegrable phase factors and global formulation of gauge fields" *Physical Review D* **12**, 3845-57