

Time-dependent symmetries: the link between gauge symmetries and indeterminism

DAVID WALLACE

1 Introduction

Mathematically, gauge theories are extraordinarily rich — so rich, in fact, that it can become all too easy to lose track of the connections between results, and become lost in a mass of beautiful theorems and properties: indeterminism, constraints, Noether identities, local and global symmetries, and so on.

One purpose of this short article is to provide some sort of a guide through the mathematics, to the conceptual core of what is actually going on. Its focus is on the Lagrangian, variational-problem description of classical mechanics, from which the link between gauge symmetry and the apparent violation of determinism is easy to understand; only towards the end will the Hamiltonian description be considered.

The other purpose is to warn against adopting too unified a perspective on gauge theories. It will be argued that the meaning of the gauge freedom in a theory like general relativity is (at least from the Lagrangian viewpoint) significantly different from its meaning in theories like electromagnetism. The Hamiltonian framework blurs this distinction, and orthodox methods of quantization obliterate it; this may, in fact, be genuine progress, but it is dangerous to be guided by mathematics into conflating two conceptually distinct notions without appreciating the physical consequences.

The price paid by this article for abandoning the mathematics of gauge theory as far as possible is an inevitable loss of rigour. Virtually nothing will be ‘proved’ below; at most, the shape of proofs will be gestured at and strong plausibility-arguments advanced. For a more detailed understanding of the mathematics, the natural place to start is Earman’s contribution to this volume (to which my own article can be seen as a commentary). Further details can be found in many standard texts on general relativity or quantum field theory (Peskin and Schroeder (1995) is particularly clear; for a really in-depth mathematical analysis, consult Henneaux and Teitelboim (1992)).

A note on terminology: the word ‘gauge’, used extensively in this introduction, will not often appear again. Its meaning is now thoroughly ambiguous (as Earman notes) and I felt it simpler to resort to the marginally more cumbersome, but clearly definable, notion of a ‘theory with a local symmetry group’. As will become clear below, there are genuine dynamical differences between

general relativity and more ‘conventional’ gauge theories such as electromagnetism, but these differences are best appreciated on their own merits rather than being annexed to the essentially sterile debate as to whether or not general relativity is ‘really’ a gauge theory. As Humpty Dumpty has taught us, words mean just what we choose them to mean — neither more nor less.

2 Symmetries of the action

The basic setup of Lagrangian mechanics is the following. We are given some system, whose configuration is specified by a point in some *configuration space*, \mathcal{Q} . (The simplest example to keep in mind is ordinary N -particle mechanics, in which the configuration space is the space of all possible sets of positions for the N particles. However, the description is just as applicable to field theories provided we are prepared to abandon the demand for a manifestly covariant description: for a field theory, a point in \mathcal{Q} is a specification of the field values at all spatial points for a given time.)¹

A path through \mathcal{Q} then specifies a possible history of the system: points on the path are labelled by time, so that the point labelled t gives the configuration of the system at time t . The only ontologically primary entities in this picture are the configurations and the paths through them: momentum, for instance, is only a derivative property of a path, and (unlike in Hamiltonian mechanics) cannot be regarded as on a par with configuration.

The task of Lagrangian mechanics is then to tell us which of the possible histories are allowed by the dynamics: that is, which histories are physically rather than merely logically possible. This is accomplished by means of the *action*, which is a rule assigning to each path γ a number $S[\gamma]$; the functional form of the action encodes everything there is to know about the dynamics of the system. The rule for specifying dynamically possible trajectories is then as follows:

1. Pick any two points (say, q_1 and q_2), and two times t_1 and t_2 .
2. Consider an arbitrary path (i. e., history) γ between q_1 and q_2 such that q_1 (q_2) is the state of the system at t_1 (t_2).
3. Evaluate the action $S[\gamma]$ for the path, as well as the actions $S[\gamma + \delta\gamma]$ for all small modifications $\delta\gamma$ of the path γ that keep the end-points q_1 and q_2 fixed.
4. γ is a dynamically possible history only if the action is extremal under variations of the form above; that is, if it is extremal in the space of histories connecting $(q_1; t_1)$ with $(q_2; t_2)$.

¹See Goldstein (1980) or any classical mechanics textbook for a discussion of Lagrangian mechanics; Goldstein (1980) also gives an introduction to infinite-dimensional systems such as fields. For more mathematically rigorous treatments of finite-dimensional Lagrangian mechanics, consult Arnold (1989) or Abraham and Marsden (1978). Introductory treatments of the mathematically rigorous theory of infinite-dimensional systems are somewhat scarce, but chapter 3 of Marsden and Ratiu (1994) is one place to start.

One naive way of defining determinism might be as follows:

A) A system is *naively deterministic* iff through any two time-labelled configurations $(q_1; t_1)$ and $(q_2; t_2)$ there exists only one dynamically possible history.

A moment's thought shows that this is too strong a restriction: consider, for instance, the dynamics of Earth's rotation around the sun, where q_1 is the location of Earth on January 1, 2000AD. Then through the time-labelled configurations $(q_1; 1/1/2000AD)$ and $(q_1; 1/1/2001AD)$ there exist (at least) two dynamically possible histories: the actual history, and one in which Earth orbited the sun in the opposite direction.

Problems like this stem from a more general difficulty in the use of the least-action principle to describe dynamics: it treats dynamics as a boundary-value problem, where histories are specified by their initial and final configurations, rather than as an initial-value problem where sufficient information about the system at only one end of the history is enough to specify the rest of the history. This suggests, however, a natural improvement to our definition of determinism:

B) A system is *manifestly deterministic* iff an arbitrarily short segment of a dynamically possible history is sufficient to fix uniquely the rest of the history.

In the rest of this article we will investigate the link between symmetry and the violation of these two forms of determinism.

3 The definition of a symmetry

For our purposes, a (variational) symmetry can be defined as a transformation \mathcal{T} on the space of possible histories, such that

1. the transformation takes each history to another history with the same action: $S[\mathcal{T}(\gamma)] = S[\gamma]$.
2. The initial and final times (t_1, t_2) of a history are left unchanged.
3. The transformation is locally defined, in the sense that the transformation \mathcal{T} for any path may be found by breaking that path into arbitrarily small components and applying \mathcal{T} separately to each component.

Obviously, the first of these is the conceptually central requirement, with the second and third being merely technical. The third condition does, however, play a crucial role in making the notion of a variational symmetry non-trivial; without it, a perfectly arbitrary permutation among all paths of a given action would count as a symmetry. (It has nothing to do with the distinction between "global" and "local" symmetries.)

In many cases the third rule can be implemented by requiring the symmetry to be defined in terms of some map f on the configuration space, so that

$\mathcal{T}(\gamma)(t) \equiv f \cdot \gamma(t)$; most familiar symmetries (such as rotation or translation, or the electromagnetic gauge symmetry) can be specified in this way. However, not all can: the diffeomorphism symmetry of general relativity cannot, for instance, and in general nor can any symmetry which corresponds to some form of time translation.

For a consideration of determinism, the most important property of a symmetry is this: since it is action-preserving, in particular it preserves the extremality of a path. In other words, if γ is a path whose action is extremised under endpoint-preserving variations, so is $\mathcal{T}(\gamma)$. But since extremality is the necessary and sufficient condition for a path to count as a dynamically possible history, it follows that $\mathcal{T}(\gamma)$ is a dynamically possible history iff γ is. That is, symmetries map the space of dynamically possible histories onto itself.

4 Global and local symmetries

For our purposes, a global symmetry group is a group (in the technical sense) of symmetry transformations which can be specified in a time-independent way: that is, the form of the transformation on a given path does not depend on the initial and final times on that path. Rotation and translation are examples of global symmetries; so is time translation in non-relativistic particle mechanics, where each particle's position is transformed to the position it would occupy Δt seconds later.

A local symmetry group, on the other hand, is a group of symmetry transformations whose action is time-dependent: that is, a group of transformations which can act independently on different segments of a path. The crucial difference between global and local symmetries is that if a global symmetry changes any segment of a history then (generically) it will change all of that history, whereas local symmetries may act non-trivially on only one segment of a history.

Our specification of a local symmetry as a time-dependent symmetry (i. e., one which is local in time) differs from the usual definition, in which the symmetry is taken as being dependent on both time and space. Bringing out the importance of *time* dependence is a major reason why this article uses the ordinary Lagrangian formalism, rather than the covariant Lagrangian formulation one more normally used in field theory.² As we will see in the next section, it is time-dependence and not spacetime dependence which is crucial to the breakdown of determinism in the presence of a local symmetry.

²By a 'covariant Lagrangian formulation' of a dynamical theory I mean one where the action S is expressed, not as an integral over time of some Lagrangian function $L(t)$, but as an integral over *spacetime* of some Lagrangian density $\mathcal{L}(\mathbf{x}, t)$. Thus the Lagrangian density of the Klein-Gordon field is $\mathcal{L}(\mathbf{x}, t) = \frac{1}{2}(\partial^\mu \phi(\mathbf{x}, t)\partial_\mu \phi(\mathbf{x}, t) - \phi(\mathbf{x}, t)^2)$, whilst the Lagrangian function $L(t)$ is the integral of $\mathcal{L}(\mathbf{x}, t)$ over all \mathbf{x} .

5 Symmetries and the breakdown of determinism

Either local or global symmetries can lead to a breakdown of naive determinism. Suppose that γ is some path from q_1 to q_2 , and suppose that \mathcal{T} is a symmetry which keeps the end-points of this particular γ unchanged: that is, $\mathcal{T}(\gamma)(t_1) = \gamma(t_1) = q_1$, and similarly for t_2 . Then, since $\mathcal{T}(\gamma)$ is dynamically possible iff γ is, and since both have the end-points, naive determinism is violated.

This phenomenon occurs only for certain global symmetries, and then only for certain initial and final points: it does not occur at all for translational symmetry, for instance, and occurs for rotational symmetry only when the initial and final points are invariant under the rotation group or one of its subgroups. (The failure of naive determinism in the case of Earth's orbit, given above, occurs because the initial and final states are both invariant under rotation about the line between Earth and the Sun.)

It is, however, ubiquitous for local symmetries: given any two end-points, we can always find elements of the local symmetry group which leave those end-points fixed but change other parts of the paths between them.

Can a similar argument show the breakdown of manifest determinism? For global symmetries, no: the time-independence of a global symmetry means that if it is trivial on the initial part of arbitrary paths through q_1 , it must be trivial on the whole path. (See figure 1). But this is not the case for local symmetries, which can perfectly well leave one segment of a path fixed and change another. (See figure 2). We can conclude, then, that whenever a theory has a local symmetry, that theory violates manifest determinism.

This is our reason for defining a local symmetry as one which is local in time rather than spacetime: local symmetries are interesting because of the breakdown of manifest determinism which they lead to, and nothing analogous occurs in spatially local symmetries. The interrelation between space and time which we have been accustomed to since Einstein and Minkowski should not blind us to the different dynamical roles which the two play.

(The purely spatial locality of a symmetry has no particularly interesting dynamical consequences: a symmetry which is spatially local but temporally global does not in any way threaten manifest determinism.³ Inevitably, in a relatively covariant theory temporal locality probably implies spatial locality and vice versa, but as long as we are interested in the dynamics the distinction between space and time remains conceptually crucial.)

³I offer two concrete examples. In electromagnetism, the group of time-independent gauge transformations is a perfectly well-behaved global symmetry, leading to an infinite set of conserved quantities $\nabla \cdot \mathbf{E}(\mathbf{x})$, one for each \mathbf{x} , but not constraining these quantities to vanish or reducing the degrees of freedom of the theory. Similarly, in general relativity the group of time-independent spatial diffeomorphisms leads to the conservation of the momentum constraints but not their vanishing.

t

(q_2, t_2)

(q_1, t_1)

\mathcal{Q}

Figure 1: A global symmetry transformation

t

(q_2, t_2)

(q_1, t_1)

\mathcal{Q}

Figure 2: A local symmetry transformation

6 Two ways to repair determinism

Earman, in his contribution to this volume, regards it as an open question as to whether or not a given theory is deterministic; we will employ a complementary strategy, imposing determinism by *fiat* and determining the consequences for a theory with a local symmetry group.

Our strategy will be as follows: let \mathcal{T}_1 be any element of the local symmetry group which leaves the initial and final parts of paths fixed and changes other parts of them. The failure of manifest determinism occurs because any dynamically possible history will be taken \mathcal{T}_1 to another such history with the same initial segment.

Our strategy for restoring determinism, then, is in essence simple:

Whenever two histories are thus related by \mathcal{T}_1 , they are in reality *the same history*.

This implies, of course, that there is not a one-to-one correspondence between the mathematics and the physics, a characteristic property of gauge theories which is discussed in Redhead's contribution to this volume.

One way of implementing this strategy — call it Option A — is to insist on the following:

Option A: Whenever two configurations are related by the action of \mathcal{T}_1 on paths through those configurations, they are in reality *the same configuration*.

This strategy is especially palatable when the symmetry itself is defined in terms of a map on the configuration space, instead of being given directly as a map on the space of paths. It is applied, for instance, in Earman's example of Maxwellian spacetime, where two configurations related by a translation are taken to be the same configuration. (It is also applicable to electromagnetism: two four-potentials related by a gauge transformation are in general taken to be the same physical state.)

It is not, however, the only strategy: consider, instead,

Option B: Two histories related by \mathcal{T}_1 furnish descriptions of the same history in terms of two different sequences of configurations, without any claim that two different points in \mathcal{Q} are the same configuration.

General relativity provides the example *par excellence* of Option B: the configurations of general relativity are 3-geometries⁴ and no two mathematically distinct 3-geometries are treated as physically the same, yet our freedom to foliate a manifold in many different ways lets us describe one and the same spacetime in terms of very different sequences of three-geometries. For a more

⁴I ignore, for simplicity, the purely spatial diffeomorphisms; if these are included then the configuration space becomes a space not of 3-geometries but of 3-metrics, and Option A applies to 3-metrics related by a purely spatial diffeomorphism.

mundane example, consider Barbour’s relational mechanics, in which time is defined only intrinsically to a path.⁵ The time labels on paths in Barbour’s configuration space, then, are arbitrary, and there exists a group of local symmetries which in effect change the time labels while keeping the points on the path fixed; thus, ‘the configuration of the system at time t ’ will have changed, but this change is due only to a redescription of the history.

There is a very substantial conceptual difference between these two ways of reinterpreting a system in the light of a local symmetry. Option A essentially means supposing that configuration space is a redundantly large description of the actual set of configurations, suggesting that the ‘real’ theory lives on some sort of quotient of configuration space in which symmetry-connected configurations are identified; such a reformulation of the theory would no longer admit either indeterminism or local symmetry. No such reformulation is available for option B: in general relativity, for example, each 3-geometry is a perfectly legitimate configuration and if some of them are purged from the configuration space of general relativity then it will not be possible to formulate the dynamics. All are needed to formulate the dynamics.

As will be seen in the next section, this conceptual difference is essentially obliterated when we move from Lagrangian to Hamiltonian mechanics, and from classical to quantum theory.

7 Local symmetries from the Hamiltonian perspective

The transformation from the Lagrangian to the Hamiltonian viewpoint is reviewed in Earman’s contribution to this volume, so I will not go into the details here. The important point about the transformation, for our purposes, is that a history of the system is now taken to be a path not in configuration space, but in phase space: that is, in the space not just of configurations but of the momenta conjugate to those configurations. As was mentioned above, this elevates momenta to the same ontological status as configuration; it is essentially equivalent to regarding the state of a system at a given time as being given not by its configuration at that time alone, but by the configuration and velocity jointly.

The local symmetry again leads to a breakdown of manifest determinism, and again we can recover determinism by insisting that mathematically different histories with the same initial and final conditions are in reality the same history. There are natural analogues of options A and B to implement this requirement:

Option A’: Whenever two states are related by a local symmetry transformation, they are in reality the same state. (This suggests that we should pass to the quotient space of phase space, dividing out by the action of the local symmetry.)

⁵For discussions of Barbour’s approach, see Barbour (1994), Pooley and Brown (2002), and references therein.

Option B’: Whenever two histories are related by a local symmetry transformation which keeps the initial and final conditions of the history unchanged, those two histories are in reality descriptions of the same history in terms of two different sequences of states.

However, option A’ is far more natural than option B’, for the following reason: specifying a history in terms of its initial and final conditions is not a natural strategy in Hamiltonian mechanics. This is because a single state, encoding as it does both position and velocity information, is in general (i. e. , when manifest determinism holds) already enough information to specify a unique dynamical trajectory. To give *two* such states, one as the initial and one as the final state, is then to over-determine the problem. If we wish to specify a history in this way we must give only the initial and final *configurations*, which breaks the symmetry between position and momentum which is so natural in Hamiltonian mechanics.

Furthermore, option A’ is less disastrous than option A for theories such as general relativity. Recall that applying option A to GR trivialises it, reducing its configuration space to a single point; however, if we are to identify only *states* and not configurations, related by the action of the symmetry, then the theory remains contentful.

However, applying option A’ to general relativity (or to Barbourian mechanics) remains conceptually unnatural. There is a perfectly natural interpretation of the local symmetry’s action on states in GR: it represents time-evolution of those states to future states of the same history, and the freedom to choose which element of the symmetry group to apply corresponds to our freedom to define time in general relativity in many ways (i. e. , to foliate spacetime in many ways). This is very different from the interpretation of the symmetry in, say, Maxwellian spacetime or electromagnetism, where it is interpreted as telling us simply that some apparently different states are really different descriptions of the very same state.

This would seem to suggest that in spite of its mathematical naturalness, we should avoid applying option A’ to theories such as general relativity. However, when we try to quantize a theory in the standard way, option A’ moves from being mathematically natural to mathematically compulsory: the Dirac quantization algorithm⁶ requires the quantum state to be invariant under the action of the local symmetry. This forces the symmetry to be interpreted as a mere redescription of the physics, rather than as a physically meaningful transformation.

We can see, then, that the standard (“canonical”) approach to the quantization of general relativity⁷ is led by the mathematics of the quantization process to interpret the diffeomorphism symmetries in a way which is conceptually quite unnatural. It is resistance to this strategy which is a prime motivation in Barbour’s alternative approach to quantum gravity.

⁶A thorough treatment of Dirac quantization may be found in Henneaux and Teitelboim (1992); see Matschull (1996) for a very clear account of the ideas involved.

⁷See Wallace (2000) for an elementary introduction to canonical quantum gravity.

8 Conclusion

The conceptually important difference between local and global symmetries is that the former, but not the latter, seem to imply a failure of determinism; this failure can in turn be traced to the fact that local symmetries allow us to transform only the mid-part of a system's history, keeping its initial and final states fixed. For this reason, it is conceptually helpful to *define* a local symmetry as one with this property — that is, as a temporally local symmetry. The spatial locality of the symmetry can be understood as conceptually uninteresting (at least from the dynamical point of view), a mere consequence of covariance.

Restoring determinism to a theory with a local symmetry requires us to drop the assumption that mathematics and physics are in one-to-one correspondence, treating mathematically distinct histories as physically the same. However, there are two very distinct ways of implementing this: for theories such as electromagnetism where the symmetry acts only on the configuration space, we regard the symmetry as telling us that certain configurations are really the same configuration, whilst for theories such as general relativity where time is not an external parameter, we regard it as telling us that the same history can be described by many different sequences of configurations.

This conceptual distinction is apparently lost when we apply standard methods of quantization to theories with local symmetries. We must await a working theory of quantum gravity before we can learn whether the loss of the distinction is an important conceptual insight into gravity and time, or simply a case of following the mathematics one step too far.

References

- Abraham, R. and J. Marsden (1978). *Foundations of Mechanics* (Second ed.). Reading, Mass.: Benjamin/Cummings.
- Arnold, V. I. (1989). *Mathematical Methods of Classical Mechanics* (Second ed.). New York: Springer. Translated by K. Vogtmann and A. Weinstein.
- Barbour, J. B. (1994). The timelessness of quantum gravity: I. The evidence from the classical theory. *Classical and Quantum Gravity* **11**, 2853–2873.
- Goldstein, H. (1980). *Classical Mechanics* (Second ed.). Reading, Mass.: Addison-Wesley.
- Henneaux, M. and C. Teitelboim (1992). *Quantization of Gauge Systems*. Princeton, NJ: Princeton University Press.
- Marsden, J. E. and T. S. Ratiu (1994). *Introduction to Mechanics and Symmetry*. New York: Springer-Verlag.
- Matschull, H.-J. (1996). Dirac's canonical quantization program. Available online from <http://www.arXiv.org/abs/quant-ph/9606031>.
- Peskin, M. E. and D. V. Schroeder (1995). *An Introduction to Quantum Field Theory*. Reading, Massachusetts: Addison-Wesley.

- Pooley, O. and H. R. Brown (2002). Relationism rehabilitated? I: Classical mechanics. *British Journal for the Philosophy of Science* **53**, 183–204.
- Wallace, D. (2000). The quantization of gravity: an introduction. Available online from <http://www.arXiv.org/abs/gr-qc/00040005>.