# Symmetry and Gauge Freedom*

Gordon Belot
Department of Philosophy
New York University
New York, NY10003
*gordon.belot@nyu.edu*

December 11, 2003

**Abstract**

The classical field theories that underlie the quantum treatments of the electromagnetic, weak, and strong forces share a peculiar feature: specifying the initial state of the field determines the evolution of some degrees of freedom of the theory while leaving the evolution of some others wholly arbitrary. This strongly suggests that some of the variables of the standard state space lack physical content—intuitively, the space of states of such a theory is of higher dimension than the corresponding space of genuine physical possibilities. The structure of such theories can helpfully be characterized in terms of the action of symmetry groups on their space of states; and the conceptual problems surrounding their strange behavior can be sharpened in light of the observation that it is usually possible to eliminate the redundant variables associated with these symmetries—which turn out to be precisely those variables whose evolution is unconstrained by the dynamical laws of the theory. This paper discusses this approach, uses it to frame questions about the interpretation of classical gauge theories, and to reflect (pessimistically) on our prospects of reaching satisfactory answers to these questions.

*Keywords:* Symmetry; Gauge Theories; Hamiltonian Mechanics; Symplectic Reduction

# 1. Introduction

Two types of theory are commonly referred to as 'gauge theories': Yang-Mills theories and constrained Hamiltonian theories. The latter class properly contains the former. Under either usage, one will observe that the most striking feature of gauge theories is that they exhibit *gauge freedom*—their initial value problems fail to be well-posed in a peculiar non-disastrous way. The differential equations of most well-behaved classical theories have well-posed initial value problems: specifying (sufficiently smooth) values for the dynamical variables at some initial time determines the values of these variables at all times; that is, there is only one solution of the equations consistent with a given set of initial data. This is *not* the case for gauge theories: the equations of such theories have infinitely many solutions for each set of initial values of their dynamical variables. In general, this is a very undesirable feature of a set of equations: for one expects a classical theory to be deterministic, in that the physical state at one time determines the physical state at all times; and this is impossible in a theory with an ill-posed initial value problem if the interpretation of the theory establishes a bijective correspondence between the space of possible physical states and the space of initial data for the equations of the theory. But the initial value problem for gauge theories fails in a very special way: there is a partition of the dynamical variables of such a theory into two classes, such that specifying initial values of the full set of dynamical variables determines the evolution of the variables of the first class, while leaving the evolution of the variables of the second class wholly arbitrary (aside from the constraints imposed by continuity and differentiability). It is thus possible to view gauge theories (of either sort) as being deterministic, so long as only the variables of the first sort are taken to have any representational import—if two dynamical states (or solutions) differ only with respect to the second sort of variable, then they  must be interpreted as representing the same physical possibility (or history).

The prominence of gauge theories in contemporary physics raises a number of conceptual and interpretative questions. What is the point of gauge freedom—why do so many theories exhibiting this curious feature play a role in our physics? What is the ontology suggested by such theories—how are we to think of degrees of freedom whose evolution is left unconstrained by the dynamics of the theory? What significance for such questions has the consideration of the quantum counterparts of our classical gauge theories?

This paper is intended to make a start on such questions. I take an oblique

approach, spending most of my time characterizing and describing a class of theories intermediate in size between the class of Yang-Mills theories and the class of constrained Hamiltonian systems—roughly, the class of constrained theories which arise when the set of conserved quantities associated with the symmetries of some Hamiltonian system are set to zero. Such theories: (i) possess the prominent structural features of paradigmatic gauge theories; (ii) allow us to examine the symmetries of Yang-Mills theories and more familiar symmetries in the same setting; and (iii) have the advantage of being relatively well understood at both the classical level and the quantum level. Thus there is reason to hope they constitute the "correct" level of generality at which to approach our problems—that this class of theories provides a perspicuous setting for a philosophical discussion of gauge freedom and related conceptual territory.[1]

The next nine sections lay out this framework: four sections which develop the general picture for classical theories are followed by three sections which situate important examples (relational classical mechanics, vacuum Maxwell theory, and vacuum Yang-Mills theory) within this framework; these are followed by two sections which touch on subtleties involving quantization and singularities. Broadly speaking, the strategy of this portion of the paper is to characterize the symmetries of classical physical theories against the background of geometric mechanics; this allows a precise and general discussion of the associated conservation laws, and the procedures of constraint and reduction which they underwrite.[2]

The final three sections of the paper involve a more direct discussion of conceptual and interpretative questions. The chief theme is that the peculiar role of classical (nonabelian) gauge theories in current physics—they are of interest only insofar as they provide insight into their quantum counterparts—means

---

[1]I believe, for instance, that this setting illuminates the analogy between paradigm examples of gauge freedom and the general covariance of general relativity. This analogy is weakest in the case of 3+1 formulations of general relativity (on this issue, see, e.g., Kuchař (1988)); it becomes stronger as one moves to progressively more covariant formulations (compare Ashtekar *et al.* (1991) and Crnkovic and Witten (1987) with Gotay *et al.* (1998)).

[2]The exposition of this material is pitched at a philosophical reader who has assimilated the basic notions of differential geometry through the geometry of connections on bundles, as presented in, say, Göckeler and Schücker (1987). For helpful philosophical discussions of this territory, see Healey (2001), Liu (2001), and Redhead (2002).

I aim for rigor within reason: given the goal of discussing finite dimensional theories and field theories in the same setting, it is crucial that definitions and results encompass the infinite dimensional case; but many details—such as functional analytic niceties—are left in the references.

that interpretative questions concerning such theories must be approached with especial care, and can seldom be expected to have determinate answers; but that, nonetheless, there is some hope that developments in the quantum theory can provide traction for interpretative questions rooted in the classical domain.

## 2. Hamiltonian Mechanics

Our starting point is the Hamiltonian formalism of classical mechanics.[3] In the simplest case, the dynamics of a set of particles is given by:

$$\text{HAMILTON'S EQUATIONS} \qquad \dot{q}^i = \frac{\partial H}{\partial p_i}, \quad \dot{p}_i = -\frac{\partial H}{\partial q_i}.$$

Here $q^i$ and $p_i$ are the positions and momenta of the particles, and the Hamiltonian, $H = H(q,p)$, is the total energy of the system in the state $(q,p) = (q^1, \ldots ; p_1, \ldots)$—typically, just the sum of the kinetic energy and the potential energy.

This formalism can be generalized in a rather straightforward manner to treat almost any classical physical theory, opening up a route to the quantization of any such theory.[4] To this end, we introduce an arbitrary manifold, $M$, as our *phase space*—this phase space is meant to parameterize the space of dynamical possibilities for the system under consideration, and corresponds to the set $\{(q,p)\}$ in particle mechanics.[5] In constructing field theories, we employ infinite dimensional manifolds—so that our phase spaces are modelled locally

---

[3]For the material sketched in this section and in §§3–5 below, see Marsden (1992), Marsden and Ratiu (1994), Schmid (1987), and the references therein.

Outside of a brief discussion in §11, the Lagrangian formalism plays no direct role in this paper. Considerations of space preclude a thorough discussion of both approaches. And it is the Hamiltonian formalism that provides the clearer perspective on our problems: the Hamiltonian version of Noether's theorem is the more straightforward and powerful; and reduction is considerably more elegant on the Hamiltonian side.

[4]In our framework, the space of states is a manifold and the dynamics are given by the flow associated with a vector field on that space. This is a familiar setting for the ordinary differential equations of classical mechanics, but it is only relatively recently that the partial differential equations of field theory and continuum mechanics have found a home there. See Olver (1993, 461–463) for discussion and references.

[5]As we will see in §10 below, it is sometimes helpful to drop even the requirement that $M$ be a manifold, and require only that the phase space be a union of nicely meshing manifolds.

upon an infinite dimensional Banach space rather than upon some $\mathbb{R}^n$.[6] As in the paradigm case of classical particle mechanics, we want our dynamics to be determined once we have specified a Hamiltonian, $H : M \to \mathbb{R}$, assigning to each dynamical state its total energy. Thus we require a means to associate with a real-valued function on $M$ a unique dynamical trajectory through each point (or, equivalently, a vector field on our phase space which can be integrated to yield such dynamical trajectories).

There are two standard ways to proceed here. Geometrically, we can observe that a Hamiltonian, $H$, has a differential, $dH$, which is, of course, a one-form on $M$. So we would be set if we had a means of associating vector fields on $M$ with one-forms on $M$: for then we could take our dynamical trajectories to be the integral curves of the vector field associated with $dH$. Our means of implementing this strategy is to introduce a closed two-form, $\omega$, on $M$, and to use this to associate a vector field, $X_f$, with any smooth function, $f$, on $M$ by solving for $X_f$ in $\omega(X_f, \cdot) = df$.[7] In order that this procedure be well-defined, we demand that $\omega$ be *non-degenerate* in that the map $v \mapsto \omega_x(v, \cdot)$, which sends vectors $v \in T_x M$ to covectors in $T_x^* M$, be injective for each $x \in M$ (so that there is at most one $X_f$ which solves $\omega(X_f, \cdot) = df$ for any $f$).[8] Such an $\omega$ is called a *symplectic form*; we call a manifold equipped with such a form a *symplectic manifold*. Our dynamics is determined by solving for the *Hamiltonian vector field, $X_H$*, in our

FUNDAMENTAL DYNAMICAL EQUATION      $\omega(X_H, \cdot) = dH,$

then integrating to find the dynamical trajectories.

Note that with $\omega$ in hand, we can define a new operation on the set of smooth functions on $M$: for $f, g \in C^\infty(M)$, let $\{f, g\} = \omega(X_f, X_g)$. This makes $C^\infty(M)$ into a Lie algebra satisfying Leibniz's rule, $\{fg, h\} = f\{g, h\} + g\{f, h\}$. We can re-write our dynamics in terms of this bracket, by declaring that $\dot{f} = \{f, H\}$ for any function, $f$, on $M$. Since the observables of our theory are represented by functions on our phase space, this completely characterizes the dynamics (which

---

[6]See Lang (1999) for a unified treatment of finite dimensional and infinite-dimensional differential geometry.

[7]We require that $\omega$ be anti-symmetric in order to implement conservation of energy (in terms of the Poisson bracket discussed below, the anti-symmetry of $\omega$ is equivalent to $\dot{H} = \{H, H\} = \omega(X_H, X_H) = 0$). We require that $\omega$ be closed to ensure that it too is preserved by dynamical evolution; see, e.g., Bates and Weinstein (1997, 23–24).

[8]We will later relax this requirement.

5

is just that given by the fundamental dynamical equation). In particular, we can plug coordinates into this equation to recover a local description of the dynamics, thus recovering Hamilton's equations (or their appropriate generalization for the system at hand).

A second approach to dynamics proceeds directly via the Poisson structure of the set of smooth functions on our phase space, without reference to the geometrical structure of that space. We begin with the set of smooth functions, $C^\infty(M)$, on some space $M$. Equipping this set with a Poisson bracket—i.e., making it into a Lie algebra obeying Leibniz's rule—determines the dynamics via $\dot{f} = \{f, H\}$. This approach is strictly more general than the first: given a Poisson bracket, $\{\cdot, \cdot\}$ on $C^\infty(M)$, it is sometimes true that $M$ is a manifold and that there is a symplectic structure, $\omega$, on $M$ such that $\{f, g\} = \omega(X_f, X_g)$ for all smooth $f$ and $g$; but there are many Poisson spaces which are not symplectic manifolds.[9]

We will call a manifold together with a Hamiltonian and either a symplectic or a Poisson structure a *Hamiltonian system*. We are primarily interested in a very special subclass of Hamiltonian systems: *simple mechanical systems*. These are specified by selecting a Riemannian manifold $(Q, g)$ together with a smooth real-valued function, $V$, on $Q$; the corresponding simple mechanical system is denoted $(T^*Q, g, V)$.[10] We will call $Q$ our *configuration space*; we take its cotangent bundle, $T^*Q$, as our phase space. We note that the metric $g$ defines a function, $T$, on $T^*Q$ via $T : (q, p) \mapsto g_q(p, p)$; this give the kinetic energy for the state $(q, p)$. Taking the function $V$ as the potential energy, we stipulate that our Hamiltonian is $H(q, p) = T(q, p) + V(q)$.

Any cotangent bundle comes equipped with a canonical symplectic structure, which we employ to determine our dynamics. Each cotangent bundle is equipped with a *canonical symplectic potential*, a one-form, $\alpha$, defined by $\alpha((q, p))(v) = p(\pi_* v)$, where $\pi : T^*Q \to Q$ is the projection $(q, p) \mapsto q$, and $\pi_*$ is the tangent map of $\pi$. The *canonical cotangent bundle symplectic form* is then $\omega = -d\alpha$; we call the associated Poisson structure the *canonical Poisson structure*. In finite dimensions these objects assume a familiar form: choosing arbitrary coordinates $\{q_i\}$ on $Q$, and writing covectors $p \in T^*_x Q$ as $p = p_i dq^i$, so that $(q^i; p_i)$ provide

---

[9]See fn. 16 below for examples where $M$ is a manifold which is not symplectic; see §10 below for cases where $M$ is not a manifold. In both sorts of case, $M$ is a union of symplectic manifolds whose Poisson brackets mesh to give the Poisson bracket on $M$.

[10]It suffices that the metric be only weakly non-degenerate, in the sense that it induces a merely injective map from $T_x Q$ to $T^*_x Q$ for each $x \in Q$; of course for finite dimensional $Q$ this will amount to the same thing as requiring a bijection.

canonical coordinates on $T^*Q$, we find that $\alpha = p_i dq^i$ and $\omega = dq^i \wedge dp_i$; the associated Poisson bracket is then just $\{f, g\} = \frac{\partial f}{\partial q^i} \frac{\partial g}{\partial p_i} - \frac{\partial g}{\partial q^i} \frac{\partial f}{\partial p_i}$; and Hamilton's equations are the equations of motion.

A diffeomorphism on a symplectic manifold which preserves the symplectic form is also called *canonical*. Coordinates in which the symplectic form or Poisson brackets assume the forms noted in the previous paragraph are also termed *canonical* (the $q$'s and $p$'s of such coordinates are said to be *canonically conjugate* to one another). Darboux's theorem tells us that every finite dimensional symplectic manifold looks locally like a cotangent bundle, in the sense that it is possible to choose a set of coordinates $(q^i; p_i)$ in which the symplectic form is just $dq^i \wedge dp_i$.[11] It follows that every finite dimensional symplectic manifold has an even number of dimensions.

## 3. Symmetries and Mechanics

Let $G$ be a Lie group, with $0 < \dim G \leq \infty$. Recall that an action of $G$ on a manifold, $M$, is a smooth homomorphism, $\Phi : G \to \text{Diff}(M)$, from $G$ to the group of diffeomorphisms of $M$; that is such an action associates each $g \in G$ with a smooth map on $M$, which we write as $\Phi_g(x)$ or, sloppily, as $g \cdot x$.[12] We call an action *proper* if the convergence of $\{x_i\}$ and $\{g_i \cdot x_i\}$, for sequences $\{x_i\} \subset M$ and $\{g_i\} \subset G$, implies the convergence of $\{g_i\}$. An action is *free* if $\Phi_e$ is the only $\Phi_g$ which fixes any points of $M$; a weaker condition is that it be *fair*—that the isotropy subgroups for any two points of $M$ are conjugate to one another in $G$.[13]

Let us suppose that we are given a simple mechanical system $(T^*Q, g, V)$ together with a Lie group, $G$ which acts on $Q$ by isometries which leave $V$ invariant; we further suppose, for reasons which will emerge by and by, that this action is both proper and fair.[14] In this case, we say that $(T^*Q, g, V, G)$ is a *simple mechanical G system*. The action of $G$ on $Q$ lifts to an action on $T^*Q$ (if we know how $g \in G$ acts on points of $Q$, then we know how it acts

---

[11]An infinite dimensional generalization of Darboux's theorem holds for those symplectic manifolds for which $v \mapsto \omega_x(v, \cdot)$ is a bijection at each $x$.

[12]Equivalently: we require $\Phi$ to be smooth and to satisfy $(gh) \cdot x = g \cdot (h \cdot x)$ and $e \cdot x = x$ for $g, h \in G$, $x \in M$, and $e$ the identity in $G$.

[13]The *isotropy subgroup* of a point $x \in M$ is $G_x = \{g \in G : g \cdot x = x\}$; subgroups $H$ and $K$ of $G$ are conjugate if there is a $g \in G$ such that $K = gHg^{-1}$.

[14]In this context, the action of $G$ proper iff $\Phi(G)$ is closed in the full group of isometries of $(Q, g)$.

on curves in $Q$, and hence how it acts on vectors and covectors). The lifted action is proper and fair. And, of course, it leaves invariant our Hamiltonian—since $H$ is the sum of two terms which are by stipulation $G$-invariant. Less obviously, the canonical cotangent symplectic potential and symplectic form are also invariant under the $G$ action.[15] Thus, the action of $G$ preserves all of the structure employed in defining our dynamics on $(T^*Q, g, V)$. It follows that $G$ maps dynamical trajectories to dynamical trajectories.

If $x, y \in T^*Q$ are such that $x = g \cdot y$ for some $g \in G$, then $x$ and $y$ are qualitatively identical from the point of view of dynamics: they share all of their dynamically relevant properties, and hence play identical roles in the structure of $(T^*Q, g, V, G)$. This implies, of course, that if we write our dynamics in terms of some coordinates around $x$ and $y$ which are related by $g$ then we find that the corresponding coordinate expressions for the dynamical trajectories around these two points assume the same form.

Taking into account the fact that $G$ is a continuous group, we can say a bit more. The orbit, $G \cdot x = \{y \in T^*Q : y = g \cdot x \text{ for some } g \in G\}$, of $x$ under the action of $G$ is a *regular* submanifold of $T^*Q$ . This means that for any $x \in T^*Q$, we can choose local coordinates of the form $\{x_1, \ldots; z_1, \ldots\}$ on some neighborhood $U \subset T^*Q$ of $x$, such that $G \cdot x \cap U = \{z_1 = 0, z_2 = 0, \ldots\}$. Thus setting the $z_i$ to zero while allowing the $x_i$ to vary carries one along the orbit of $x$. But each point of $G \cdot x$ is qualitatively identical. Because the points of such orbits are dynamically indifferent, the $x_i$ are dynamically irrelevant—any way of setting their values leads to the "same" evolution. This suggests in turn that it may be possible to drop the $x_i$ from our theory altogether.

This would amount to identifying $G$-related points of $T^*Q$, and projecting down to the resulting quotient space the Hamiltonian and the canonical Poisson structure. Something like this is indeed possible, and the space which results encodes the $G$-invariant dynamics of our simple mechanical $G$ system. But this space—a Poisson manifold which is not symplectic—has a rather convoluted structure and lies off of our present track.[16] So let us set that strategy aside,

---

[15]Indeed, the lift of any diffeomorphism of $Q$ to $T^*Q$ preserves the symplectic potential, and hence also our symplectic form. In the finite dimensional case, this follows from the fact that the expression for the symplectic potential assumes the same form in every coordinate system on $Q$.

[16]Here is an instructive sort of example; for details see Arnold (1989, Appendix 2) or Marsden and Abraham (1970).

The configuration space for a rigid body constrained to rotate about a fixed point is just $SO(3)$; the phase space is $T^*SO(3)$. The Hamiltonian is just the kinetic energy induced by a (right)

8

and pursue a less ambitious one, which requires us to restrict our attention to a submanifold of $T^*Q$ before identifying $G$-related points. To this end, we next consider the conserved quantities associated with the action of $G$ on $T^*Q$.

## 4. Symmetry and Constraint

The infinitesimal generator of an element of a Lie group $G$ is an element of $\mathfrak{g}$, the Lie algebra of $G$; the infinitesimal generator of a diffeomorphism of a manifold $M$ is a vector field on $M$; so the action of a Lie group on a manifold, which associates a diffeomorphism, $\Phi_g$, of $M$ with each $g \in G$, yields a means of associating a vector field, $\xi_M$, on $M$ with each $\xi \in \mathfrak{g}$.[17]

In our case, the group $G$ acts on both the configuration space, $Q$, and the phase space, $T^*Q$. So each $\xi \in \mathfrak{g}$ is associated with both a vector field $\xi_Q$ on $Q$ and a vector field $\xi_{T^*Q}$ on $T^*Q$. In the case of $\xi_{T^*Q}$, we can ask whether it is the Hamiltonian vector field of any function on $T^*Q$—that is, for $\xi \in \mathfrak{g}$ we can look for a function $J^\xi \in C^\infty(T^*Q)$, such that $\xi_{T^*Q} = X_{J^\xi}$. We do not have to look far: $J^\xi(q, p) = \langle p, \xi_Q(q) \rangle$ does the trick (the bracket on the right hand side is the pairing between tangent and cotangent vectors at $q \in Q$).[18] The correspondence is linear. Moreover, $J^{[\xi,\zeta]} = \{J^\xi, J^\zeta\}$ so $\xi \in \mathfrak{g} \mapsto J^\xi \in C^\infty(T^*Q)$ is a Lie algebra

---

invariant metric on the group. Allowing the group to act on itself, we have a six dimensional simple mechanical $SO(3)$ system. Identifying points related to each other by the action of $SO(3)$ on the phase space yields a system living on $\mathbb{R}^3$ carrying an interesting Poisson bracket (this reduced phase space is the dual of the Lie algebra of $SO(3)$). Being odd-dimensional, $\mathbb{R}^3$ cannot be a symplectic manifold. But it *is* foliated by symplectic manifolds—the spheres of radius $R$ about the origin. The volume form on each sphere is a symplectic form. And the sphere of radius $R$ is the image under the quotient map of the set of points in $T^*SO(3)$ with angular momentum of magnitude $R$. The conservation of the magnitude of angular momentum, a consequence of the *dynamics* of the original theory is now built into the *geometry* of the reduced phase space (the Poisson structure keeps you on the sphere you start on, no matter what Hamiltonian is imposed). The resulting dynamical equations are Euler's equations for a rigid body.

A directly analogous route leads from the treatment of a perfect fluid in a given spatial domain, $\Omega$, in which one keeps track of the location of the fluid particles (so that the configuration space is $S\,\mathrm{Diff}(\Omega)$, the infinite dimensional (ILH) group of volume-preserving diffeomorphisms on $\Omega$) to familiar treatment in which one keeps track only of the velocity field of the fluid (so that the phase space is the space of divergence free vector fields on $\Omega$—i.e., the dual of the Lie algebra of $S\,\mathrm{Diff}(\Omega)$) and in which the equations of motion are Euler's equations for a perfect fluid.

[17]That is: $\xi_M(x) = \frac{d}{dt}(\exp(t\xi) \cdot x)\,|_{t=0}$ for $\xi \in \mathfrak{g}$.

[18]Note that it is more common to work with the dual *momentum map*, $\mathbf{J} : T^*Q \to \mathfrak{g}^*$ defined by $\langle \mathbf{J}(x), \xi \rangle = J^\xi(x)$ (here the pairing is between elements of the Lie algebra and of its dual).

homomorphism.[19] In the context of a simple mechanical $G$ system, each $\xi \in \mathfrak{g}$, $J^\xi$ is a conserved quantity of the dynamics: that $\dot{J}^\xi = \{J^\xi, H\} = 0$ is just the infinitesimal restatement of the invariance of $H$ under the action of $G$.

We will see below that the standard examples of conserved quantities in physics arise out of this construction.

Let $G_x$ be the isotropy subgroup of some arbitrary $x \in T^*Q$, consisting of the elements of $G$ which fix $x$ (so that $G_x = \{e\}$ for all $x \in T^*Q$ if the action of $G$ on $T^*Q$ is free). Then the fact that $G$ acts fairly implies that $\dim G_x$ is independent of $x$. It turns out that $\mathfrak{g}$ is mapped by $J$ onto a subspace of $C^\infty(T^*M)$ of dimension $\dim_{\text{act}} G := \dim G - \dim G_x$.[20] Thus we have a strong Hamiltonian version of Noether's theorem: choosing a basis $\{\xi_i\}$ for $\mathfrak{g}$ gives us $\dim_{\text{act}} G$ independent conserved quantities, which generate a Poisson algebra that is a homomorphic image of $\mathfrak{g}$. When, as in the examples of §§6–8 below, the action is free, we obtain $\dim G$ conserved quantities, whose Poisson algebra is isomorphic to $\mathfrak{g}$.

Our strategy is to start with a simple mechanical $G$ system $(T^*Q, g, V, G)$ , then to investigate the dynamics which arises when we restrict attention to the constraint surface $\Gamma = \{x \in T^*Q : 0 = J(\xi_1)(x) = J(\xi_2)(x) = \ldots\}$, with $\{\xi_i\}$ a basis for $\mathfrak{g}$.[21] That is, we set the conserved quantities associated with the group action to zero.[22]

There are two perspectives from which we can study the dynamics induced on $\Gamma$ by our fundamental dynamical equation, $\omega(X_H, \cdot) = dH$. Thinking of $\Gamma$ as a subset of $T^*Q$—which latter we henceforth call the *extended phase space*—, we know that the symplectic form and Hamiltonian of the ambient space determine a unique trajectory through each point of $T^*Q$; of course, trajectories through

---

[19]Something stronger is true. $J$ is *equivariant* with respect to the adjoint action of $G$ on $\mathfrak{g}$: $J^{\mathrm{Ad}_g \xi}(g \cdot x) = J^\xi(x)$ (recall that for matrix groups, $\mathrm{Ad}_g \xi = g\xi g^{-1}$). Equivalently, the momentum map, $\mathbf{J}$, is equivariant with respect to the coadjoint action of $G$ on $\mathfrak{g}^*$: $\mathrm{Ad}^*_{g^{-1}} \circ \mathbf{J} = \mathbf{J} \circ \Phi_g$. Much of the theory of conserved quantities and reduction sketched below can be recovered in more general contexts—such as proper group actions on arbitrary symplectic or Poisson manifolds—so long as the momentum maps involved are equivariant.

[20]That is, $\dim_{\text{act}} G$ is just the dimension of any orbit of the action. Here and below, expressions involving algebraic operations on dimensions should be interpreted informally when infinite dimensional objects are in play. In §10 below we will be interested in a case where $\dim G$ is infinite while $\dim G_x$ is finite; in this case, as our (now ill-formed) formula suggests, we recover infinitely many conserved quantities.

[21]$\Gamma$ is independent of the basis chosen. Equivalently, $\Gamma = \mathbf{J}^{-1}(0)$.

[22]In each of the applications of §§6–8 below, there is strong physical motivation for imposing the constraint.

points of $\Gamma$ lie entirely on $\Gamma$, since the constraint surface is characterized by the vanishing of a set of conserved quantities.

It turns out that this perspective involves treating $\Gamma$ *extrinsically*—it makes essential use of information about the behavior of the Hamiltonian and the symplectic structure off of $\Gamma$. This becomes clear when we consider the alternative, *intrinsic* approach to the dynamics on $\Gamma$. Let us consider $\Gamma$ as a manifold in its own right, and demand that the objects appearing in our fundamental dynamical equation be defined on $\Gamma$: we require $X_H$ to take values in $T\Gamma$, and replace $H$ by $H \mid_\Gamma$, by restricting the arguments of the Hamiltonian function to points on the constraint surface. We accomplish this by noting that the embedding $i : \Gamma \hookrightarrow T^*Q$ gives us the ability to identify vectors in $T\Gamma$ with vectors tangent to $\Gamma$ in $T^*Q$, and allows us to pull back $\omega$ to a form, $\omega \mid_\Gamma := i^*\omega$, defined on $\Gamma$.[23] We are interested in vector fields, $X_H$, living on $\Gamma$ that solve $\omega \mid_\Gamma (X_H, \cdot) = d(H \mid_\Gamma)$.

Now, $\omega$ is a symplectic form on $T^*Q$: a closed, non-degenerate two-form. It follows immediately that $\omega \mid_\Gamma$ is a closed two-form on $\Gamma$. But $\omega \mid_\Gamma$ is degenerate: at each point $x \in \Gamma$, the set of null vectors, $\{v \in T_x\Gamma : \omega(v, w) = 0 \quad \forall w \in T_x\Gamma\}$, is just the set of infinitesimal generators of the $G$ action at that point, $\{\xi_{T^*Q}(x) : \xi \in \mathfrak{g}\}$. Thus the null vectors are the tangent vectors to the orbits of the action of $G$.

The dynamical trajectories on $\Gamma$ determined by the extrinsic dynamical problem give us a set of curves on $\Gamma$. These curves will be integral curves of some vector field $X_H$ tangent to $\Gamma$ that solves the intrinsic dynamical problem. But $X_H$ is by no means the only solution to the intrinsic version of the fundamental dynamical equation. The others arise as follows. Let $N(x)$ be a null vector field on $\Gamma$. Then $X_H + N$ also solves the intrinsic dynamical equation:

$$\omega \mid_\Gamma (X_H + N, \cdot) = \omega \mid_\Gamma (X_H, \cdot) + \omega \mid_\Gamma (N, \cdot)$$

$$= \omega \mid_\Gamma (X_H, \cdot)$$

$$= d(H \mid_\Gamma).$$

Thus at each point of $\Gamma$, the dynamical equation determines the tangent vector to the dynamical trajectory only up to the addition of an arbitrary infinitesimal

---

[23]The tangent map $i_* : T\Gamma \to T^*Q$ is an isomorphism on to its image. For $v, w \in T\Gamma$, $\omega \mid_\Gamma (v, w) = \omega(i_*v, i_*w)$.

generator of the action of $G$. Upon integration, this means that dynamical trajectories determined by the intrinsic equation are of the form $h(t) \cdot x(t)$, where $x(t)$ comes from a solution of the extrinsic dynamical equation, and $h : \mathbb{R} \to G_0$ is a smooth function which selects for each time, $t$, an element of the connected component of the identity in $G$.[24] In particular, since we can always choose $h(0) = e$ (where $e$ is the identity in $G$), we see that there is a $\dim_{\text{act}} G$ dimensional family of intrinsic dynamical trajectories through each point $x$ of $\Gamma$ (each of which, in fact, arises as a solution of the extrinsic dynamical equation, for some Hamiltonian, $H^*$, on $T^*Q$ that restricts to $H$ on $\Gamma$).

So the initial value problem for the dynamics on $\Gamma$, considered intrinsically, is well-posed only up to a time-dependent transformation from $G_0$—and specifying a point $x$ in $\Gamma$ determines the evolution of only those quantities which are invariant under the action of $G_0$.

## 5. Symmetry and Reduction

We now identify points on $\Gamma$ related by the action of $G$.[25] If $x, y \in \Gamma$ with $y = g \cdot x$ for some $g \in G$, then we write $x \sim y$. This is an equivalence relation; we denote the equivalence class of $x$ by $[x]$. We are interested in the quotient space, $\Gamma/G = \{[x] : x \in \Gamma\}$ (equipped with the quotient topology whose open sets are those with open pre-images in $\Gamma$).

Let $x, y \in \Gamma$ with $x = g \cdot y$; and let $x(t)$ and $y(t)$ be the extrinsic dynamical trajectories through $x$ and $y$. Since the action of $G$ maps dynamical trajectories to dynamical trajectories, it follows that $x(t) = g \cdot y(t)$ for each $t \in \mathbb{R}$—i.e., $[x] = [y]$ implies $[x(t)] = [y(t)]$. So extrinsic dynamical trajectories through points of $\Gamma$ which project to the same point of $\Gamma/G$ themselves project to the same curve in $\Gamma/G$.

Now consider an intrinsic dynamical trajectory on $\Gamma$. This will be of the form $h(t) \cdot x(t)$, with $h : \mathbb{R} \to G_0$, for some extrinsic dynamical trajectory, $x(t)$. But $[h(t) \cdot x(t)] = [x(t)]$ for each $t$, so $x(t)$ and $h(t) \cdot x(t)$ project down to the *same* curve in $\Gamma/G$. Indeed, by the same reasoning the complete pencil of intrinsic

---

[24]Recall that only elements in $G_0$ can be reached by exponentiating elements of $\mathfrak{g}$.

[25]The construction described in this section is known as *cotangent bundle reduction*. Like its generalizations, *symplectic reduction* (alias *Marsden-Weinstein reduction*) and *Poisson reduction*, this procedure has its roots in $19^{th}$ century techniques for eliminating variables from dynamical problems, but has only been developed in a fully general and global form since the 1970s.

dynamical trajectories through a given point of $\Gamma$ projects down to a single curve in $\Gamma/G$.

So we have a unique curve through each point of $\Gamma/G$, the image of all of the dynamical trajectories, intrinsic as well as extrinsic, through all of the points on $\Gamma$ which project down to that point of $\Gamma/G$. Remarkably, these curves on $\Gamma/G$ are generated by a Poisson structure and Hamiltonian which $\Gamma/G$ inherits from the embedding of $\Gamma$ in $T^*Q$. We take the smooth functions on $\Gamma/G$ to be the restrictions to $\Gamma$ of the $G$-invariant functions on $T^*Q$. This latter set has a Poisson structure, which we can now think of as the Poisson structure of $C^\infty(\Gamma/G)$.[26] Since our original Hamiltonian on $T^*Q$, $H$, is $G$-invariant, it projects down to a well-defined function, $\tilde{H}$, on $\Gamma/G$. As usual, our Poisson structure and Hamiltonian determine a set of dynamical trajectories—which are just the images of the dynamical trajectories of $\Gamma$ under the projection to $\Gamma/G$.[27]

In the present context—where our dynamical system has a nice cotangent bundle structure and our group acts properly and fairly on $Q$—we can say quite a bit more about $\Gamma/G$ and the dynamics defined upon it. Consider the quotient space $Q/G$ which results when we identify points of $Q$ related by the action of $G$. We call this space the *reduced configuration space*; it is a manifold of dimension $\dim Q - \dim_{\text{act}} G$.[28] Because the metric, $g$, and potential, $V$, on $Q$ are $G$-invariant, they project down to yield a metric, $\tilde{g}$, and potential, $\tilde{V}$, on $Q/G$. Thus we have a simple mechanical system $(T^*(Q/G), \tilde{g}, \tilde{V})$. We call the cotangent bundle $T^*(Q/G)$ the *reduced phase space*. Remarkably, the dynamical theory which we have reached via quotienting $Q$ by the action of $G$ coincides with the one constructed above by quotienting $\Gamma$ by the action of $G$ (i.e., there is an isomorphism which preserves all relevant structures).

Proceeding via either route, we end up with a dynamical system living on

---

[26]This is induced from the canonical Poisson structure of $C^\infty(T^*Q)$ by taking the quotient of the algebra $G$-invariant functions on the extended phase space by the ideal of functions which vanish on the constraint surface.

[27]The picture developed in this paragraph holds for a very general class of systems—one needs a Hamiltonian system on a symplectic or Poisson manifold, together with a proper group action which admits an equivariant momentum map. See Bates and Lerman (1997) and Ortega and Ratiu (1998).

[28]Note that the dimension of the orbits of the actions of $G$ on $Q$ and $T^*Q$ are equal. When $\dim Q = \dim_{\text{act}} G = \infty$, $\dim Q - \dim_{\text{act}} G$, interpreted informally, can turn out to be either infinite or finite. Indeed, Yang-Mills theories provide examples of both sorts of behavior: see §§7 and 8 below for infinite dimensional reduced configuration spaces; see Gotay (1989) and Rajeev and Rossi (1995) for finite dimensional ones.

a phase space of dimension $2(\dim Q - \dim_{\text{act}} G)$ which encodes the $G$-invariant dynamics of the sector of the original theory in which the conserved quantities associated with the action of $G$ vanish.

We now have three spaces in play, each carrying its own dynamics: the extended phase space, $T^*Q$; the constraint surface, $\Gamma$; and the reduced phase space, $\Gamma/G \simeq T^*(Q/G)$. We are interested in two quotients: the quotient of $\Gamma$ by $G$ is the reduced phase space, the quotient of $Q$ by $G$ is the reduced configuration space. Note, finally, that $\Gamma \to \Gamma/G$ and $Q \to Q/G$ are principal $G$ -bundles when the action of $G$ is free.[29]

## 6. Relational Mechanics

We can now turn to a concrete example. The idea is to examine the action of the group of Euclidean isometries on a system of $N$ gravitating Newtonian point particles.[30] Here identifying points in the configuration space related by the action of the group means taking as the space of possible configurations the space of possible relative distances between the particles. The corresponding phase space is the space of possible relative distances and relative velocities of the particles—so the reduced theory is, in an interesting sense, a fully relational dynamical theory.[31] Of course, the discussion of the previous section shows that this theory can also be viewed as the result of restricting attention to dynamical states for which the conserved quantities (i.e., the linear and angular momenta of the system) vanish, then identifying points related by isometries, and projecting down the Hamiltonian and the Poisson structure.

In slightly more detail, we begin by considering $\mathbb{R}^{3N}$, the space of possible dispositions of the $N$ particles in three dimensional Euclidean space. This

---

[29]Or, rather, this is known to hold in the finite dimensional case. Some extra work is required for infinite dimensional systems; see Mitter (1980) for the Yang-Mills case. Note further that because $Q$ carries a metric, the bundle $Q \to Q/G$ carries a natural connection; for applications and discussion, see Batterman (2002), Koiller *et al.* (1996) and Littlejohn and Reinsch (1997).

[30]Galileian boosts can be handled in generalizations of the present framework, but do not satisfy all of our present conditions. (i) They do not arise as lifts of transformations of the configuration space. (ii) While they leave invariant the dynamical trajectories of the theory, they do not leave the Hamiltonian itself invariant.

[31]I take the liberty of speaking of velocities rather than momenta. The sort of theory discussed here can be constructed via other methods. For a distinctive and influential approach see Barbour and Bertotti (1982); for an insightful commentary on that scheme, see Pooley and Brown (2002).

space carries a flat metric, $g$, encoding the kinetic energy of the particles, and a potential function, $V$, encoding the gravitational potential energy. Thus we have a simple mechanical system, $(T^*\mathbb{R}^{3N}, g, V)$. We are interested in the action of $E(3)$, the group of Euclidean isometries (i.e., products of shifts, rotations, and reflections). This group acts in an obvious way on the configuration space, $\mathbb{R}^{3N}$—simply shifting, rotating, or reflecting the position of each of the particles in physical space. Of course, this action leaves invariant the kinetic energy (which derives from the Euclidean metric on the space in which the particles move) and the potential energy (which sees only relative distances between particles), and lifts to a canonical action on $T^*\mathbb{R}^{3N}$ (translations leave the momenta of each particle invariant, while rotations and reflections act upon them in the obvious way).

But the action of E(3) on $\mathbb{R}^{3N}$, while proper, is not fair. Generic points of $\mathbb{R}^{3N}$ are not fixed by any non-trivial isometries; but there is a set of measure zero consisting of points which have higher symmetry (e.g., szygial configurations—in which the particles are collinear—are fixed by rotations about the line on which the particles lie). In order to obtain a fair (indeed, free) action, we restrict, until §10 below, our attention to $Q$, the set of generic points of $\mathbb{R}^{3N}$ representing asymmetric particle configurations.[32] Of course, $g$ and $V$ still live on $Q$, so we can take $(T^*Q, g, V, E(3))$ as our simple mechanical $G$ system.

The conserved quantities associated with the Euclidean symmetry of our theory are just (the components of) the linear momentum and the (center of mass) angular momentum of the system. So our constraint surface, $\Gamma$, comprises the states with vanishing linear and angular momentum. From the extrinsic point of view, the dynamical initial value problem on $\Gamma$ is, of course, well-posed—while from the intrinsic point of view, the motion of the system is determined only up to a time-dependent translation and rotation (the connected component of the identity, $SE(3)$, is generated by translations and rotations).[33] Equivalently, at each point on $\Gamma$, the null space of the restriction of the symplectic form of $T^*Q$ to $\Gamma$ is spanned by the infinitesimal generators of rotations and translations.

Identifying points on $\Gamma$ related by Euclidean isometries leads to a dynamical theory on $T^*(Q/E(3))$ which captures the $E(3)$ -invariant dynamics of systems

---

[32] We ought to also exclude collision singularities. That leaves us with a Hamiltonian vector field that is smooth but incomplete. This is unpalatable, but appears to be unavoidable, since it is known that some collision singularities are unregularizable; see Abraham and Marsden (1978, 699) or Diacu and Holmes (1996, Chapter 3).

[33] See Lynden-Bell (1995) for a Lagrangian formulation corresponding to our Hamiltonian account of the intrinsic dynamics on $\Gamma$.

of particles with vanishing linear and angular momentum. The manifold $Q/E(3)$ is the space of spatial configurations of the particles modulo isometries. A point in this space is specified by the set of relative distances between the particles. Now, the number of relative distances is $N(N-1)/2$ while $Q/E(3)$ is $3N-6$ dimensional. For $N >> 4$, the former number is much larger than the latter, and the relative distances provide a vastly *over-complete* set of coordinates on $Q/E(3)$. Correspondingly, $T^*(Q/E(3))$ is (in general) over-completely parameterized by the set of relative distances and velocities of the particles.

So the fact that our reduced theory has a well-posed initial value problem means that specifying the initial values of the relative distances and velocities determines their past and future values—something which is not true in general in Newtonian mechanics, of course, but is true when the total angular momentum and linear momentum vanish. Since, our cosmos does, in fact, appear to have vanishing angular momentum, this relational theory enjoys the same degree of empirical adequacy as Newtonian celestial mechanics.[34]

# 7. Vacuum Maxwell Theory

We now consider the simplest Yang-Mills theory: vacuum electromagnetism.[35] We forget about relativity—so our theory describes the evolution in time of the electric and magnetic field on physical space, $S$ rather than the behavior of the electromagnetic field on spacetime. For convenience, we assume that $S$ is a flat Riemannian three manifold (typically, a three-torus or $\mathbb{R}^3$). In this setting, Maxwell's equations for the magnetic and electric fields, $\mathbf{B}$ and $\mathbf{E}$, are just:

$$
\begin{aligned}
(i) && \dot{\mathbf{B}} &= -\operatorname{curl}\mathbf{E} \\
(ii) && \dot{\mathbf{E}} &= \operatorname{curl}\mathbf{B} \\
(iii) && \operatorname{div}\mathbf{B} &= 0 \\
(iv) && \operatorname{div}\mathbf{E} &= 0.
\end{aligned}
$$

These equations can be derived within our framework.[36] A vector potential, $\mathbf{A}$, is a smooth vector field, $\mathbf{A} : S \to \mathbb{R}^3$, on physical space. We take as our

---

[34]The orthodox view is that the cosmic background radiation puts a fantastically small bound on the magnitude of the angular momentum of the universe; see Obukhov (2000) for references, orthodox and otherwise.

[35]See Marsden and Weinstein (1982), Schmid (1987), or Marsden and Ratiu (1994) for details.

[36]It may bear emphasizing that it is possible to formulate a Hamiltonian theory of the behavior

configuration space $\mathcal{A} = \{\mathbf{A}\}$, the space of vector potentials on $S$. $\mathcal{A}$ is an infinite dimensional vector space.[37] Because $\mathcal{A}$ is a vector space, we can identify $T_{\mathbf{A}}\mathcal{A}$, the tangent space at $\mathbf{A}$, with $\mathcal{A}$ itself. Thus a cotangent vector at $\mathbf{A}$ is something which eats elements of $\mathcal{A}$ and spits out real numbers—for instance, another vector field $\mathbf{E} : S \to \mathbb{R}^3$ with the pairing between vectors and covectors given by integrating the scalar product of $\mathbf{A}$ and $\mathbf{E}$ over $S$. So we can take as our phase space $T^*\mathcal{A} = \{(\mathbf{A}, \mathbf{E}) : \mathbf{A}, \mathbf{E} : S \to \mathbb{R}^3\}$, itself an infinite dimensional vector space.

$T^*\mathcal{A}$ carries the canonical cotangent bundle symplectic and Poisson structures. These are given by obvious analogs of the finite dimensional formulae:

$$\omega((\mathbf{A}_1, \mathbf{E}_1), (\mathbf{A}_2, \mathbf{E}_2)) = \int_S (\mathbf{E}_1 \cdot \mathbf{A}_2 - \mathbf{E}_2 \cdot \mathbf{A}_1) dx$$

and

$$\{F, G\} = \int_S \left( \frac{\delta F}{\delta \mathbf{E}} \frac{\delta G}{\delta \mathbf{A}} - \frac{\delta F}{\delta \mathbf{A}} \frac{\delta G}{\delta \mathbf{E}} \right) dx.$$

(See Schmid (1987) or Marsden and Ratiu (1994) for functional derivatives such as $\frac{\delta F}{\delta \mathbf{E}}$.) We take $H(\mathbf{A}, \quad \mathbf{E}) = \frac{1}{2} \int_S |\mathbf{E}|^2 + |\text{curl}\,\mathbf{A}|^2 \, dx$ as our Hamiltonian. Note that this is, as usual, the sum of a kinetic term (arising from a flat metric on $\mathcal{A}$) and a potential term (given by a scalar on $\mathcal{A}$)—so we have in hand a simple mechanical system.

Hamilton's equations on $T^*\mathcal{A}$ are just $\dot{\mathbf{A}} = -\mathbf{E}$ and $\dot{\mathbf{E}} = \text{curl}\,\text{curl}\,\mathbf{A}$. Defining, as usual, $\mathbf{B} := \text{curl}\,\mathbf{A}$, we see that these equations of motion imply the first two Maxwell equations. The third follows from the identity $\text{div}\,\text{curl}\cdot = 0$. The fourth, $\text{div}\,\mathbf{E} = 0$, will emerge shortly.

Our next step is to consider the action of the group, $\mathcal{G}$, of gauge transformations on our theory. This group is just (a subgroup, picked out by appropriate boundary conditions, of) the additive group of smooth functions on

---

of $\mathbf{B}$ and $\mathbf{E}$ directly, avoiding the excursion via the vector potential. That this is possible in practice is an interesting and important feature of Maxwell theory. That it is possible in principle follows from the discussion of §5.

[37]Here and below all infinite dimensional spaces are taken to be appropriate Sobolev spaces; see references. The scalar potential plays no role in this approach—from the perspective of the standard route to Maxwell's theory, this amounts to choosing to work in the temporal gauge (i.e., we fix the gauge; see fn. 79 below). But no conceptual questions are begged, since freedom to perform (spatial) gauge transformations quickly re-emerges.

$S$, and acts on $\mathcal{A}$ via $f : \mathbf{A} \mapsto \mathbf{A} + \operatorname{grad} f$.[38] The corresponding action on $T^*\mathcal{A}$ is $f : (\mathbf{A}, \mathbf{E}) \mapsto (\mathbf{A} + \operatorname{grad} f, \mathbf{E})$. Being a lift, this action preserves the canonical cotangent bundle symplectic form on $T^*\mathcal{A}$. Leaving $\mathbf{E}$ invariant it leaves the kinetic energy invariant. It also leaves invariant the potential energy: $|\operatorname{curl}(\mathbf{A} + \operatorname{grad} f)|^2 = |\operatorname{curl} \mathbf{A}|^2$, because curl grad $\cdot = 0$. But $\mathcal{G}$ fails to act fairly on $\mathcal{A}$ and $T^*\mathcal{A}$—symmetric fields have larger-than-generic isotropy groups. In order to circumvent this problem, we restrict our attention to the subgroup of $\mathcal{G}$ consisting of *pointed gauge transformations*: $\mathcal{G}_* = \{f \in \mathcal{G} : f(x_0) = 0\}$ for some arbitrary but fixed $x_0 \in S$.[39] $\mathcal{G}_*$ acts fairly (indeed, freely) on $\mathcal{A}$ and $T^*\mathcal{A}$.

Thus we have a simple mechanical $\mathcal{G}_*$ system. Since our group is infinite dimensional, there are infinitely many conserved quantities associated with our symmetries. These are encoded in the function div $\mathbf{E}(x)$ which remains constant at each point $x \in S$ as $\mathbf{E}$ evolves in accord with the equations of motion.[40]

In accord with our usual procedure, we now set these conserved quantities to zero, and investigate the constraint surface $\Gamma = \{(\mathbf{A}, \mathbf{E}) \in T^*\mathcal{A} : \operatorname{div} \mathbf{E} = 0\}$. Thus, restricting attention to $\Gamma$ amounts to imposing the fourth Maxwell equation! The null directions of the restriction of the symplectic form of $T^*\mathcal{A}$ to $\Gamma$ correspond, of course, to the infinitesimal generators of the action of $\mathcal{G}_*$ on $T^*\mathcal{A}$. So the intrinsic equations of motion on $\Gamma$ are well-posed only up to a time-dependent (small) gauge transformation: if $(\mathbf{A}(t), \mathbf{E}(t))$ is a solution then so is each $(\mathbf{A}(t) + \operatorname{grad} g(t), \mathbf{E}(t))$, for $g : \mathbb{R} \to \mathcal{G}_{*_0}$.[41] As always, the extrinsic dynamics has a well-posed initial value problem.

Identifying gauge-related points on $\Gamma$ leads to a Hamiltonian theory on the reduced phase space: $T^*(\mathcal{A}/\mathcal{G}_*)$ equipped with its canonical cotangent bundle symplectic structure, and carrying the Hamiltonian $H(\mathbf{A}, \mathbf{E}) = \frac{1}{2} \int_S |\mathbf{E}|^2 + |\mathbf{B}|^2 \, dx$. Once again, the Hamiltonian of the reduced theory is the sum of a kinetic term arising from a metric on the reduced configuration space and of a potential term defined on the reduced configuration space. Thus we have recovered a simple mechanical system as our reduced theory.

The reduced configuration space $\mathcal{A}/\mathcal{G}_*$ is the infinite dimensional manifold

[38]Here and below all groups of gauge transformations are Hilbert Lie groups; see Schmid (1987).

[39]Note that we lose little in the shift from $\mathcal{G}$ to $\mathcal{G}_*$: $\mathcal{G}_*$ is a normal subgroup of $\mathcal{G}$, and $\mathcal{G}/\mathcal{G}_*$ is just the one dimensional group $U(1)$.

[40]For a statement in line with our official notion of conserved quantities in terms of $J^\xi$, see Śniatycki (2000, equation 5).

[41]Gauge transformations lying in the connected component of the identity are referred to as *small gauge transformations*.

of vector potentials modulo pointed gauge transformations. In the special case where $S = \mathbb{R}^3$, this space is just the space of possible magnetic fields (i.e., divergence free $\mathbf{B} : S \to \mathbb{R}^3$); and the reduced phase space can be taken to be $\{(\mathbf{B}, \mathbf{E}) : \operatorname{div} \mathbf{B} = \operatorname{div} \mathbf{E} = 0\}$, carrying the Poisson bracket

$$\{F, G\} = \int_S \left( \frac{\delta F}{\delta \mathbf{E}} \operatorname{curl} \frac{\delta G}{\delta \mathbf{B}} - \frac{\delta G}{\delta \mathbf{E}} \operatorname{curl} \frac{\delta F}{\delta \mathbf{B}} \right) dx.$$

(Note that $\mathbf{B}$ and $\mathbf{E}$ are not canonically conjugate with respect to the cotangent bundle Poisson structure on $T^*(\mathcal{A}/\mathcal{G}_*)$.) The equations of motion are then just $\dot{\mathbf{B}} = -\operatorname{curl} \mathbf{E}$ and $\dot{\mathbf{E}} = \operatorname{curl} \mathbf{B}$—we recover the Maxwell equations in their most familiar form.

If, however, space is non-simply connected (e.g., a torus), then specifying the magnetic field fails to determine a point in $\mathcal{A}/\mathcal{G}_*$: there are vector potentials $\mathbf{A}$ and $\mathbf{A}'$ such that curl $\mathbf{A}$=curl $\mathbf{A}'$ but there is no $f \in \mathcal{G}_*$ such that $\mathbf{A}' = \mathbf{A}+\operatorname{grad} f$—so specifying a magnetic field fails to determine a gauge-equivalence class of vector potentials. In this case, a better parameterization of $\mathcal{A}/\mathcal{G}_*$ is provided by the set of *holonomies*: for each closed curve $\gamma$ starting and ending at our basepoint $x_0 \in S$, we define the holonomy of $A$ around $\gamma$, $H_\gamma(\mathbf{A}) = \exp i \oint_\gamma \mathbf{A} \, ds$. For any $f \in \mathcal{G}_*$, $H_\gamma(\mathbf{A}) = H_\gamma(\mathbf{A} + \operatorname{grad} f)$, so $H_\gamma$ is a gauge-invariant quantity on $\mathcal{A}$ for each $\gamma$. In fact, two vector potentials yield the same holonomies iff they are gauge-equivalent. So the set of holonomies provides a good set of coordinates on $\mathcal{A}/\mathcal{G}_*$. But this set, like the set of relative distances in relational mechanics, forms a vastly over-complete set of coordinates on the reduced configuration space. Just as relative distances must satisfy, e.g., the triangle inequality, so must the holonomies satisfy certain constraints (these are surprisingly elegant; see Barrett (1991)).

Of course, one might well prefer to avoid spacetime non-local quantities such as holonomies. And this may appear to be possible in the case at hand: even when the magnetic fields fail to exhaust the content of the reduced configuration space, it remains true that the familiar Maxwell equations for $\mathbf{B}$ and $\mathbf{E}$ exhaust the content of the Hamiltonian equations of motion. Surely, then we are justified in eschewing the representational resources of the reduced configuration space, insofar as they outrun those afforded by the space of magnetic fields? But trouble is not far away: the magnetic field alone fails to contain all of the information necessary to construct an account of a quantum charged particle moving in a classical electromagnetic field on a non-simply connected space.[42]

---

[42]Because of the Aharonov-Bohm effect; see Belot 1998 for this story.

And, in any case, there is no avoiding such non-local quantities once one moves from Maxwell's theory to non-abelian Yang-Mills theories.

## 8. Vacuum Yang-Mills Theory

We now turn to more general vacuum Yang-Mills theories. Let $M$ be a compact three dimensional Riemannian manifold representing physical space. Let $G$ be a compact connected finite dimensional Lie group; from the compactness of $G$, it follows that we can equip the Lie algebra, $\mathfrak{g}$, with an inner product, $\langle \cdot, \cdot \rangle$, invariant under the adjoint action of $G$ on $\mathfrak{g}$; for notational convenience, we assume that $G$ is a matrix group, so that $\mathrm{Ad}_g \xi = g\xi g^{-1}$. Let $P \to M$ be a principal $G$-bundle over $M$. The configuration space for our Yang-Mills theory will be the space of connections on $P$ —physically, the space of field potentials.[43]

We focus, initially, on the case where $P$ is a trivial bundle (the machinery necessary for non-trivial $P$ is sketched below, set off from the main text by square brackets). In this case we can fix an arbitrary trivialization of $P$; this allows us to pull back forms on $P$ to forms on $M$, so that we can take the relevant fields to live on physical space.

Recall that, just as an ordinary $p$-form on $M$ eats vector fields on $M$ and yields a real number for each $x \in M$, so a $\mathfrak{g}$ -valued $p$-form on $M$ eats vector fields and yields an element of $\mathfrak{g}$ for each $x \in M$. The wedge product of a real-valued one-form and a real-valued two-form on $M$ is a real-valued three-form—which, because $M$ is three dimensional, we can think of as being the multiple by some real-valued function of the volume form on $M$ associated with the Riemannian metric on $M$. We want a similar mechanism for pairing $\mathfrak{g}$ -valued one- and two-forms on $M$ to produce a real-valued function on $M$. We achieve this by employing the following rule: let $\alpha$ and $\beta$ be, respectively, a $\mathfrak{g}$-valued one-form and a $\mathfrak{g}$-valued two-form on $M$; and let $\{\xi_i\}$ be a basis for $\mathfrak{g}$; then we can write $\alpha = \sum_i \alpha_i(x) \otimes \xi_i$ and $\beta = \sum_i \beta_i(x) \otimes \xi_i$, where each $\alpha_i$ is a real-valued one-form on $M$, and each $\beta_i$ is a real-valued two-form on $M$; employing the inner product on $\mathfrak{g}$, we can now construct the pairing $\langle\langle \alpha, \beta \rangle\rangle = \sum_{i,j} \langle \xi_i, \xi_j \rangle\, \alpha_i \wedge \beta_j$; this yields the desired real-valued three-form on $M$. Note, in particular, that this allows

---

[43]For this approach to Yang-Mills theories, see Arms (1981), Mitter (1980), and Moncrief (1980). Śniatycki (1999) and Śniatycki *et al.* (1996) provide details for the case where $P \to M$ is trivial. Maxwell theory is the special case where $G = U(1)$. The complications introduced by taking sources into account are briefly touched upon in §11 below. We are again choosing the temporal gauge; see fn. 37 above.

us to construct a norm, $\|\alpha\| = \int_M \sqrt{\langle\langle \alpha, *\alpha \rangle\rangle}$, on the spaces of $\mathfrak{g}$-valued one- and two-forms over $M$ (here $*\alpha$ is $\mathfrak{g}$-valued $(3-p)$-form dual to $\alpha$ by the Hodge star).

In our present context, we can think of connections on $P$ as $\mathfrak{g}$-valued one-forms on $M$. We take the space of such connection one-forms, $\mathcal{A}$, as our configuration space. $\mathcal{A}$ is an infinite dimensional vector space; so we identify the tangent space at $A \in \mathcal{A}$ with $\mathcal{A}$ itself. Then a cotangent vector at $A \in \mathcal{A}$ eats elements of $\mathcal{A}$ and returns real numbers. So we can take a cotangent vector, $E \in T_A^*\mathcal{A}$ to be a $\mathfrak{g}$-valued two-form on $M$, with the pairing given by $\int_M \langle\langle A, E \rangle\rangle$. So our phase space is $T^*\mathcal{A} = \{(A, E)\}$, which, as usual, we equip with its canonical cotangent bundle symplectic structure.[44] As our Hamiltonian, we take $H(A, E) = \frac{1}{2} \left( \|E\|^2 + \|F_A\|^2 \right)$, where $F_A = D_A A$ is the $\mathfrak{g}$-valued curvature two-form, associated with the connection $A$ via the action of $D_A$, the covariant exterior derivative associated with $A$. Note that this Hamiltonian is, as usual, the sum of a kinetic term (deriving from a flat Riemannian metric on $\mathcal{A}$) and a potential term (given by a function on $\mathcal{A}$). So we have a simple mechanical system, with equations of motion $\dot{A} = E$ and $\dot{E} = -\operatorname{curl} B - [[A\times, B]]$ (where $B := \operatorname{curl} A + [[A\times, A]]$ and where $[[\alpha, \beta]] := \sum_{i,j}[\xi_i, \xi_j]\,\alpha_i \wedge \beta_j$ for $\alpha = \sum_i \alpha_i(x) \otimes \xi_i$ and $\beta = \sum_i \beta_i(x) \otimes \xi_i$).

We now turn to the symmetries of this theory. Depending on the geometry of $M$, the theory may be invariant under a group of spatial isometries. Let us set those aside. Our interest lies in the action of $\mathcal{G}$, the infinite dimensional group of gauge transformations. In our present context, a gauge transformation is of the form $x \in M \mapsto g(x) \in G$, a smooth assignment of an element of $G$ to each point of $M$. These act on $\mathcal{A}$ via $A(x) \mapsto g(x)^{-1}A(x)g(x) + g(x)^{-1}dg(x)$. This action lifts to an action on $T^*\mathcal{A}$ which preserves the symplectic structure; this action is just $(A, E) \mapsto (g^{-1}Ag + g^{-1}dg, g^{-1}Eg)$ (omitting this time the arguments). $F_A$, like $E$, transforms under a gauge transformation by conjugation. The invariance under gauge transformations of $\|E(x)\|$ and $\|F_A(x)\|$ then follows immediately from the invariance under conjugation of the inner product on the Lie algebra of $G$. Thus $\mathcal{G}$ is a symmetry group for our theory, leaving invariant both the symplectic structure and the Hamiltonian.

We do not quite have a simple mechanical $\mathcal{G}$ system though, for the action of $\mathcal{G}$ on $\mathcal{A}$ is proper but not fair—since symmetric field potentials have larger-than-generic isotropy groups. To get around this problem, we restrict our attention to the group, $\mathcal{G}_*$, of pointed gauge transformations, which acts freely on $\mathcal{A}$ (i.e., we

---

[44]See Śniatycki (2000, equation 3) for the symplectic potential.

require our gauge transformations to assign the identity element in $G$ to some fixed but arbitrary $x_0 \in M$). $\mathcal{G}_*$ is an infinite dimensional normal subgroup of $\mathcal{G}$ (indeed, $\mathcal{G}_*$ is a normal subgroup of $\mathcal{G}$, with $\mathcal{G}/\mathcal{G}_* = G$—so not much is lost).

The conserved quantities associated with this group action are encoded in the $\mathfrak{g}$-valued three-form $D_A E$, which we think of as a $\mathfrak{g}$ -valued function on $M$, preserved by the dynamics of the theory.[45] So our constraint surface, $\Gamma$, is given by imposing Gauss' Law: $D_A E(x) = 0$ at each $x \in M$. The null directions on this constraint surface correspond to the infinitesimal generators of gauge transformations, which are just given by assignments $x \in M \mapsto \xi(x) \in \mathfrak{g}$. We can, of course, study the dynamics on $\Gamma$ either extrinsically or intrinsically. The former point of view yields a well-posed initial value problem, while the latter yields the more familiar formulation in which the initial value problem is well-posed only up to a time-dependent (small) gauge transformation.[46]

[Before discussing the structure of the reduced configuration space, let me remark that the picture is much the same when we work in the more general context in which $P \to M$ is allowed to be a non-trivial bundle.[47] $A$, $E$, and $F_A$ must now be thought of as $\mathfrak{g}$-valued forms on $P$, while $\mathcal{G}$, the group of gauge transformations, is now the set of vertical automorphisms of $P$.

In the case where $P$ was trivial, we were able to take one or more $\mathfrak{g}$-valued forms, $\alpha, \dots$, living on $M$, perform some operation $T(\alpha, \dots)$ on them to yield a scalar function on $M$, then work with the integral $\int_M T(\alpha, \dots)$. This played a crucial role in our definition of the pairing between tangent and cotangent vectors to our phase space, and in the definition of the kinetic and potential energies. In the case at hand, we still need to integrate scalars over $M$ to define the necessary structures—but now the forms we start with are forms on $P$. We deploy a somewhat klunky apparatus to get around this problem. We choose a set of distinguished disjoint open sets $U_i \subset M$, the union of whose closures cover $M$, and over which we are able to define local sections of $P$. Given a form, $\alpha$, on $P$, we

---

[45]See Śniatycki (2000, equation 5) for the official statement.

[46]Gauge transformations lying in the connected component of the identity are again referred to as *small gauge transformations*. The fact that $\mathcal{G}$ and $\mathcal{G}_*$ are not in general connected has important consequences for quantum gauge theories; see Jackiw (1984) on $\theta$ angles.

[47]Some accounts of gauge theories leave the impression that a nontrivial $P$ is required iff a monopole is being described. This is not so. On the one hand, the 't Hooft-Polyakov monopole lives on a trivial bundle; see Göckeler and Schücker (1987, §§10.3 and 10.6). On the other hand: if our physical space is $\mathbb{R}^4$, we may treat $\infty$ as a point, compactifying $\mathbb{R}^4$ to $S^4$; in order to give a consistent global description, we may require a nontrivial $P \to S^4$, although the original bundle over $\mathbb{R}^4$ was necessarily trivial; see Singer (1982, 38–40).

can use our local sections to construct pullbacks, $\alpha_i$ living on the $U_i$; we can then perform some operation, $T(\alpha_i, \ldots)$ on our pull backs, yielding a scalar on each $U_i$; finally, we can define $\int_M T(\alpha, \ldots) = \sum_i \int_{U_i} T(\alpha_i, \ldots)$. This construction will be independent of our choice of $U_i$ and local sections, so long as each of the $T(\alpha_i, \ldots)$ is gauge invariant in the sense that $T(\alpha_i(x), \ldots) = T(g(x)\alpha_i(x)g^{-1}(x), \ldots)$ for any map $x \mapsto g(x)$ defined on $U_i$.[48]

Our configuration space, $\mathcal{A}$, will be the infinite dimensional space of connections on $P$. A connection, $A$, is a $\mathfrak{g}$-valued one-form on $P$ which has two special features: (i) $A(x)(\xi_P(x)) = \xi$ for each $\xi \in \mathfrak{g}$ and $x \in P$ (where $\xi_P$ is the vector field on $P$ associated to $\xi$ in virtue of the action of $G$ on each fibre of $P$); (ii) it is *equivariant* in that $\Phi_g^* A = \mathrm{Ad}_{g^{-1}} \circ A$ (where $\Phi_g$ is the action of $g \in G$ on $P$ and $\Phi_g^*$ is its tangent map).

$\mathcal{A}$ is an affine subspace of the space of $\mathfrak{g}$-valued one-forms on $P$. Indeed, choosing an arbitrary connection $A_0$, we can write $\mathcal{A} = \{A_0 + \alpha\}$ where $\alpha$ ranges over *tensorial* $\mathfrak{g}$-valued one-form on $P$—where a $\mathfrak{g}$-valued $p$-form on $P$ is tensorial if it is equivariant and *horizontal*, in that $\alpha(\xi_P^1, \ldots, \xi_P^p) = 0$ whenever $\xi^1, \ldots, \xi^p \in \mathfrak{g}$. For each $p$, the space of tensorial $\mathfrak{g}$-valued $p$-forms on $P$ is a vector subspace of the space of $\mathfrak{g}$-valued $p$-forms on $P$. Tensorial $\mathfrak{g}$-valued forms on $P$ have the following gorgeous property: if $\alpha$ is tensorial, and $\alpha'$ and $\alpha''$ are pullbacks to $U_i$ associated with two local sections of $P$ over $U_i$ related by the transformation $x \in U_i \mapsto g(x) \in G$, then $\alpha'(x) = g(x)\alpha''(x)g(x)^{-1}$.

Our discussion of the affine nature of $\mathcal{A}$ shows that the tangent space to $\mathcal{A}$ is just the space of tensorial $\mathfrak{g}$-valued one-forms on $P$. We can take a cotangent vector $E \in T_A^* \mathcal{A}$ at an element of $A$ to be a $\mathfrak{g}$-valued tensorial two-form on $P$; the pairing between tangent vectors and cotangent vectors is given by pulling them back to $M$ via our arbitrary local trivializations, to yield a $\mathfrak{g}$-valued one-form and a $\mathfrak{g}$-valued two-form defined on each of our privileged neighborhoods of $M$, proceeding with these pullbacks as in the case where $P$ is trivial, then summing. The construction is kosher because the pullbacks transform by conjugation, and the inner product on $\mathfrak{g}$ is invariant under conjugation.

---

[48]There is an elegant alternative to this procedure. Associated with our $P \to M$, there is a vector bundle $\mathrm{ad}\, P \to M$ with typical fibre $\mathfrak{g}$, on which $G$ acts via $\mathrm{Ad}$, and which has the same transition functions as $P$. There is a natural correspondence between the tensorial $\mathfrak{g}$-valued forms on $P$ (discussed below) and forms on $M$ which take values in the space of sections of $\mathrm{ad}\, P$ —our local pullbacks can be sewn together to yield such sections, and the necessary pairings, norms, etc. can be defined in a manifestly invariant fashion on the space of such sections. See, e.g.,Marathe and Martucci (1992, §6.4).

As usual, we equip $T^*\mathcal{A}$ with its canonical cotangent bundle symplectic structure. We construct our Hamiltonian on $T^*\mathcal{A}$ by pulling back $E$ and $F_A$ to $M$, employing the expression from the trivial case on each of our privileged neighborhoods, then summing; again it follows from the fact that $F_A$, like $E$, is tensorial on $P$ that the resulting expression does not depend on our choice of local sections. Thus we have a simple mechanical system on $T^* \mathcal{A}$.

The expressions for the transformation of $A$ and $E$ under gauge transformation are unchanged; in order to achieve a fair (indeed, free) action, we again restrict to $\mathcal{G}_*$, the group of pointed gauge transformations (required to reduce to the identity on the fiber over some fixed but arbitrary $x_0 \in M$). The conserved quantities are again encoded in the three-form, $D_A E(x)$, defined now on $P$. Imposing Gauss' law, $D_A E(x) = 0$ gives us a constraint surface $\Gamma$. The null directions on the constraint surface are given by the infinitesimal generators of the gauge transformations (given, now, by $G$-equivariant functions from $P$ to $\mathfrak{g}$). As usual, the extrinsic dynamics on this constraint surface has a well-posed initial value problem, while the initial value problem for the dynamics under the standard, intrinsic, construal is well-posed only up a time-dependent small gauge transformation.]

Whichever context we choose to work in, the infinite dimensional reduced configuration space—the space of connections modulo pointed gauge transformations—is a complicated object. For non-abelian $G$, $\mathcal{A}/\mathcal{G}_*$ is non-linear and carries a non-flat Riemannian metric.

Holonomies provide the best known parameterization of $\mathcal{A}/\mathcal{G}_*$. If we fix a point $b$ in the fibre of $P$ above our basepoint $x_0 \in M$ and a connection $A \in \mathcal{A}$, then we can associate a $g \in G$ to each smooth curve, $\gamma$, in $M$ which begins and ends at $x_0$; we say that $g$ is the holonomy of $A$ around $\gamma$ for $b$.[49] Keeping $b$ fixed, we can observe that a connection gives us a map from the space of closed curves through $x_0$ to $G$, called the holonomy map. It turns out that two connections yield the same holonomy map if and only if they differ by a pointed gauge transformation. Again, though, this set of coordinates on the reduced configuration space is highly over-complete.

We can also start from the other end: with only a little care, we can turn the space of closed curves through our basepoint into a topological group, $\mathcal{L}M_*$ (with concatenation as the multiplication; see Barrett 1991). Then any homomorphism from $\mathcal{L}M_*$ to $G$ which is appropriately smooth determines a $\mathcal{G}_*$-

---

[49]Here is the recipe: construct the horizontal lift of $\gamma$ which begins at $b$; the end point of this curve lies in the same fibre as $b$; $g$ is the element of $G$ which maps $b$ to this point.

equivalence class of connections on some principal fibre bundle $P \to M$. Thus such a smooth homomorphism determines the topology of $P$ as well as the geometry of the connection.

## 9. Quantization

Our reduced theories are standard Hamiltonian theories living on cotangent bundles, and can in principle be quantized via any procedure adapted to such theories. In practice, however, the theories discussed in §§6–8 above exhibit features which complicate the execution of standard recipes (here I am thinking of the over-completeness of the natural sets of coordinates on the reduced configuration spaces, and, in the case of Yang-Mills theories, of the non-locality of the holonomies). In such situations, an alternative approach pioneered by Dirac (1964) becomes very attractive. We first construct a quantization for the extended phase space. Then we restrict attention to a subset of states which obey a quantum version of the classical constraints—if the classical constraints are of the form $C_i(x) = 0$ for some functions, $C_i$, on the extended phase space, then the quantum constraints are $\hat{C}_i\psi = 0$ (i.e., we restrict our attention to states annihilated by the operator corresponding to $C$). The hope is that this procedure yields the same result as would a direct quantization of the reduced theory.

This hope is borne out for a large class of finite dimensional systems.[50]

Suppose that we have a finite dimensional simple mechanical $G$ system $(T^*Q, g, V, G)$, in which $G$ is connected and unimodular and the action of $G$ on $Q$ is free.[51] We undertake to construct a quantization of this system—i.e., a representation of an interesting Poisson subalgebra of $C^\infty(T^*Q)$ as an algebra of self-adjoint operators on $L^2(Q)$. In general, there is considerable choice available here, even in the finite dimensional regime to which we are presently restricting our attention—the case of $\mathbb{R}^{2n}$, where the Stone-von Neumann theorem assures

---

[50]What follows is a sketch of the central result of Gotay (1986a). See Sjamaar (1996) and Huebschmann (2002) for related results in other finite dimensional cases. See Śniatycki (2000) for the status of this program in the Yang-Mills case—where, of course, there is the additional problem of the construction of the correct inner product on the space of states. There is some reason to think that the Dirac program falters in the context of singular quotients; see Römer (1988).

[51]$G$ is *unimodular* if it carries a bi-invariant measure; for this, it suffices that $G$ be abelian or compact.

us that there is only a single irreducible representation of the usual algebra of $p$'s and $q$'s, is quite exceptional (on this point, see, e.g., Isham (1983)). We fix our attention on geometric quantization (see Woodhouse (1980) or Bates and Weinstein (1997)), where the imposition of further geometric structure upon $T^*Q$ determines a quantization. It is natural to impose the further condition that the quantization respect the structure of our theory as a simple mechanical $G$ system: that the new structure imposed be compatible with the cotangent bundle structure of $T^*Q$ and with the action of $G$.

With a quantization of the extended phase space in hand, we can proceed to construct a Dirac quantization of the constrained theory. Recall from §4 that the basis elements, $\xi_i$, for the Lie algebra $\mathfrak{g}$ of $G$ are associated with functions, $J^{\xi_i}$, on the extended phase space which generate the action of $G$ on the this space, and whose Poisson brackets represent the Lie bracket relations between the $\xi_i$. Our constraint surface, $\Gamma$, is determined by setting the $J^{\xi_i}$ equal to zero. Now, $G$ also acts on $L^2(Q)$, the Hilbert space of our quantization of the extended theory. The infinitesimal generators of that action are just the $\hat{J}^{\xi_i}$, the operators corresponding to the $J^{\xi_i}$ (this follows from the compatibility assumed above). We construct a quantization of the constrained system by imposing the *Dirac condition*: we restrict our attention to those $\psi \in L^2(Q)$ such that $\hat{J}^{\xi_i}\psi = 0$.[52] This amounts to restricting our attention to the $G$-invariant states in $L^2(Q)$.[53] The states satisfying this condition form a Hilbert space which inherits a representation of some $G$-invariant observables from the extended quantization.

Now, because we take our quantization of the extended theory to be determined by new structures which mesh with the cotangent bundle structure and the group action of the extended theory, we find that these new structures

---

[52]Complications, such as rigged Hilbert spaces, may be necessary to handle the case where 0 lies in the spectra of these operators, but not in their discrete part.

[53]The $\hat{J}^{\xi_i}$ are the infinitesimal generators of the action of $G$, so we can recover the action of $G$ by exponentiating their action (since $G$ is connected): for the states we are interested in, the $\hat{J}^{\xi_i}$ act like the zero operator; so their exponentials act like the identity.

If we lift the requirement that $G$ be unimodular, then the Dirac condition must be amended to read $\hat{J}^\xi\psi = -\frac{i}{2}\operatorname{tr}(\operatorname{ad}\xi)\psi$ (for unimodular groups, $\operatorname{tr}(\operatorname{ad}\xi) = 0$), in order to ensure that the Dirac quantization matches the standard quantization of the reduced theory (see Duval *et al.* (1990) for an example where this is essential). This formula, which ought to look surprising in light of the first consideration adduced in this footnote, arises naturally within BRST quantization; see Loll (1992) and Tuynman (1992) for accessible treatments. This scheme, with its ghost and anti-ghost variables offers a nice illustration of the theme, flagged in §13 below, that mathematical tractability can sometimes be secured by introducing non-physical degrees of freedom.

project down to $T^*(Q/G)$, inducing a quantization of the reduced theory (i.e., a representation of a Poisson subalgebra of $C^\infty(T^*(Q/G))$ as an algebra of observables on $L^2(Q/G)$). This Hilbert space is canonically unitarily equivalent to that of the corresponding Dirac quantization (the states of the constrained quantization are, after all, $G$-invariant functions in $L^2(Q)$, and thus may be thought of as functions in $L^2(Q/G)$); this isomorphism also relates the well-behaved observables of the two quantizations.[54]

This sort of result shows that one can often avoid working directly with the reduced theory. But note that it shows only that quantizations of the extended theory that respect its structure as a simple mechanical $G$ system are associated with quantizations of the reduced theory via the Dirac procedure. It does *not* follow that imposing the Dirac condition gives us a way of moving from arbitrary quantizations of the extended theory to quantizations of the reduced theory. Indeed, there exist constrained systems lying just outside of our ambit that admit geometric quantizations that: (i) fail to respect the constraint structure; and (ii) do not arise as quantizations of the reduced theory.[55]

## 10. Symmetry and Singularities

We now want to dispense with the assumption that $G$ acts fairly on $Q$. We will see that this means countenancing reduced spaces $Q/G$ and $\Gamma/G$ which are singular spaces, not manifolds.

We begin with some generalities about quotients of finite dimensional manifolds by proper group actions.[56] Let $G$ act properly on a finite dimensional manifold, $X$. Then the topological space $X/G$ (equipped with the projection topology) is Hausdorff. We count a function on $X/G$ as smooth it corresponds to a $G$-invariant smooth function on $X$. Now, $X/G$ will *not* be a manifold unless the action is fair. But it can *always* be taken to be composed of manifolds, in the following sense. Construct an equivalence relation on the set of subgroups of $G$ by declaring $H$ and $K$ to be equivalent if they are conjugate in $G$—i.e., if

---

[54]Problems may arise at this stage with the quantum Hamiltonian. This is not unusual in geometric quantization.

[55]See Ashtekar and Horowitz (1986) and Gotay (1986b) for such systems; the constraints in question are not associated with momentum maps. Loll (1992) broaches some related questions.

[56]For finite dimensional systems; see, e.g., Pflaum (2001a, Chapter 4; 2001b, §5). Much of the picture is known to carry over to interesting infinite dimensional cases; see Isenberg and Marsden (1982), Kondracki and Rogulski (1986), and Śniatycki *et al.* (1996).

there is a $g \in G$ with $K = gHg^{-1}$. Writing the equivalence class of $H$ as $(H)$, we define $X_{(H)} := \{x \in X : (G_x) = (H)\}$, the set of points of $X$ whose isotropy subgroup is in $(H)$. Each $X_{(H)}$ is a submanifold of $X$, called the *stratum of points of symmetry type* $(H)$, and $X$ is the disjoint union of the $X_{(H)}$.[57] Similarly, each $(X/G)_{(H)} := X_{(H)}/G \subset X/G$ is a manifold, the stratum of points of type $(H)$, and $X/G$ is the disjoint union the strata $(X/G)_{(H)}$.[58] We will say that $X/G$ has a *non-trivial stratum structure* if it has more than one non-empty stratum. If we consider two subgroups of $G$, $H$ and $K$, such that there is a $g \in G$ with $H \subset gKg^{-1}$, then we can say that $(H)$ corresponds to points with less symmetry than those corresponding to $(K)$. In this case, $X_{(K)}$ is contained in the boundary of $X_{(H)}$ in $X$ and $(X/G)_{(K)}$ is in the boundary of $(X/G)_{(H)}$ in $X/G$. There will be a minimal symmetry type, $(H_{\min})$, often corresponding to the identity subgroup of $G$; the corresponding sets of generic points with minimal symmetry, $X_{\mathrm{reg}} := X_{(H_{\min})}$ and $(X/G)_{\mathrm{reg}} := (X/G)_{(H_{\min})}$, are open and dense in $X$ and $X/G$, respectively.

When $X$ has the additional structure of a configuration space or a phase space, the quotients that we are interested in will inherit some corresponding structure.[59] If $Q$ is a Riemannian manifold, there is a standard technique for making $Q/G$ into a metric space: the distance between two points of the quotient space is given by taking the infimum over the lengths of curves joining the corresponding orbits in $Q$ (the minimizing curves are geodesics orthogonal to the orbits of $G$). Of course, $Q/G$ will not be a Riemannian manifold in the ordinary sense if it is singular. But it can be equipped with a sort of generalized metric tensor (Pflaum (2001a, §2.4)), and we will, in any case, be able to pursue dynamics on $Q/G$ in the form of a geodesic principle.[60]

---

[57]Strictly speaking, $X_{(H)}$ and the other strata discussed in this section may be $\Sigma$-manifolds— finite or countable disjoint unions of manifolds which need not all be of the same dimension (see Pflaum (2001a,b)). I describe the situation for the case where each stratum consists of a single connected component; the more general theory can be recovered via some fairly obvious modifications—such as relativizing the claims of the next note to connected components.

[58]Each $X_{(H)} \to (X/G)_{(H)}$ is a $G/H$ bundle, so we have $\dim(X/G)_{(H)} = \dim X - \dim G + \dim H$.

[59]For configuration spaces, see Alekseevsky *et al.* (2001) and Pflaum (2001a, §2.4) for the finite dimensional case and Kondracki and Rogulski (1986) for the Yang-Mills case. For phase spaces, see Bates and Lerman (1997) and Ortega and Ratiu (1998) for the finite dimensional case, Isenberg and Marsden (1982) for general relativity, and Śniatycki *et al.* (1996) for Yang-Mills.

[60]This is the proper setting for the theory of Barbour and Bertotti (1982), which is based upon Jacobi's principle. Note that in quantization of simple mechanical systems, the Laplacian associated with the metric on configuration space plays a prominent role. Presumably, in order

If $G$ acts properly on a configuration space, $Q$, then it also acts properly on the corresponding phase space, $T^*Q$, and constraint surface, $\Gamma$.[61] In this case each stratum, $(T^*Q)_{(H)}$, of $T^*Q$ is a symplectic submanifold of $T^*Q$; each stratum, $\Gamma_{(H)}$, is a submanifold of $T^*Q$ with a null space of dimension $\dim G - \dim H$; and each stratum, $(\Gamma/G)_{(H)}$, of the reduced phase space is a symplectic manifold. Each such quotient stratum forms a dynamically closed subset: a dynamical trajectory through a point of a given stratum is confined to that stratum (since it is a symplectic manifold). The corresponding Poisson brackets on the quotient strata mesh to give us a Poisson bracket on $\Gamma/G$. Together with the projection of the $G$-invariant Hamiltonian, this allows us to formulate a dynamical theory on the reduced phase space, $\Gamma/G$.

How does dropping the requirement that the action be fair bear upon our examples?

One consequence in the Yang-Mills case is that we can work with the full group of gauge transformations rather than with the subgroup of pointed gauge transformations. We then get $\dim G$ further conserved quantities (see, e.g., Landsman (1998a, §IV.3.6)). In the abelian case, the holonomies remain invariant under the larger group; but in the non-Abelian case, this is not so, and one has to work instead with the traces of the holonomies, the so-called Wilson loops (these are maps from loops to the complex numbers). In any case, one advantage of shifting to the larger group is that we no longer have to fuss with base-points—we can now work with the set of *all* smooth loops on $M$.[62]

There has been a good deal of work on the structure of the reduced configuration spaces of Yang-Mills theories (that is, the spaces of connections modulo the full group of gauge transformations; see, e.g., Kondracki and Rogulski (1986), Huebschmann (1996) and Rudolph *et al.* (2002)). And there has been some speculation that the singularity structure of such spaces may play an important role in non-perturbative aspects of QCD (see Asorey (1999) and Viela Mendes (2000)). There has been rather less work on the stratum structure of the reduced *phase space* of Yang-Mills theory.[63] The problems are somewhat different,

---

to quantize a singular system, one will need to find an operator related in an appropriate way to the generalized metric on the reduced configuration space.

[61]If there is more than one non-empty stratum of $T^*Q$ corresponding to an isotropy group of positive dimension, then $\Gamma$ itself possesses singularities and fails to be a submanifold of $T^*Q$.

[62]However, the characterization of those maps from the space of loops into the complex numbers which correspond to connections on fibre bundles is slightly cumbersome and depends upon $G$; see Loll (1994).

[63]See Śniatycki *et al.* (1996) for the case where physical space is three dimensional, Rajeev

as the singularities in the reduced configuration space derive from symmetric vector potentials, while the singularities in the reduced phase space derive from states in which both the vector potential and canonically conjugate electric field are symmetric.

This important difference, which we may roughly characterize by saying that the reduced phase space is not the cotangent bundle of the reduced configuration space, arises whenever the group action on our configuration space is not fair. It is easiest to visualize in the case of particle mechanics. There, allowing non-free group actions means that it is no longer necessary to excise symmetric points from the phase space in setting up the theory. In our previous treatment of this example, we took as our phase space $T^*(\mathbb{R}^{3N}_{\text{reg}})$, the set of states whose *configuration* variables were symmetry free. We now want to investigate the full phase space, $T^*\mathbb{R}^{3N}$.[64] Thus our extended phase space will include states whose configuration variables have nontrivial symmetries. So our reduced phase space will include non-empty strata corresponding to such states, each such stratum forming a dynamically closed subset. Now, one might have hoped that, as in the regular case, the reduced phase space, $(T^*\mathbb{R}^{3N})/E(3)$, would be isomorphic to $T^*(\mathbb{R}^{3N}/E(3))$, the cotangent bundle of the reduced configuration space.[65] But this is not so: the former space includes states in which the configuration variables have a higher degree of symmetry than the momentum variables; the latter does not.[66] Indeed, consider a state in which the particles are collinear, but this configurational symmetry is fleeting because the momenta exhibit no symmetry. Such states exist in $(T^*\mathbb{R}^{3N})/E(3)$, lying in the strata of generic points (since the asymmetry of the momentum variables means that the state as a whole is not left invariant under any symmetries). But these states are not to be found in $T^*(\mathbb{R}^{3N}/E(3))$: since they involve a szygial configuration, they would have to sit in the cotangent bundle of the stratum of $\mathbb{R}^{3N}/E(3)$ consisting of szygial configurations; but this cotangent bundle only has as many momentum variables as the reduced stratum of szygial configurations has configuration degrees of freedom, which means that it contains only states in which the symmetry of the momentum variables is at least as great as those of the configuration variables.

and Rossi (1995) for the case where physical space is taken to be a circle.

[64]More properly: the space which results upon the excision of points corresponding to collisions.

[65]We interpret the latter space as the disjoint union of the cotangent bundles of the strata of the reduced configuration space. For details, see Pflaum (2001a, §2.3).

[66]Both spaces include states in which the momenta exhibit a higher degree of symmetry than the configuration variables—e.g., generic configurations in which each particle is at rest.

Now, it might well be thought that this sort of disparity between the reduced phase space and the cotangent bundle of the reduced configuration space could not be of any genuine importance. The real interest of the classical theories under discussion lies, it is natural to think, in their quantum analogs. And, after all, the set of generic points of minimal symmetry is open and dense in the reduced phase space; indeed, the cotangent bundle of the quotient of the generic stratum of the reduced configuration space is open and dense in the reduced phase space. So surely the singular points, being of Borel measure zero, can make little difference to the quantum theory? This suggests that it suffices to quantize the theory on $T^*(Q_{\text{reg}}/G)$.[67]

But it is not clear that this approach is the only one, or that it is in general to be preferred to an attempt to quantize the full reduced theory, singular strata and all. Indeed, in some examples it appears that the singularity structure of the reduced phase space plays a role in determining crucial boundary conditions for the operators of the quantum theory.[68] The quantization of singular spaces remains ill-understood.[69]

## 11. A Puzzle and Two Solutions

I turn at last to direct consideration of conceptual and interpretative questions. Let us begin with *a puzzle*. We can study systems of gravitating point particles with or without imposing the (observationally underwritten) constraint that the angular and linear momenta should vanish. We can study the Yang-Mills dynamics of connections on principal bundles with or without restricting our attention to the vacuum case by imposing Gauss' Law, $D_A E = 0$, as a constraint.[70]

---

[67]See Emmrich and Römer (1990) for this approach, and for indications that quantum states have some tendency to cluster near singularities. One advantage of this approach is that $Q_{\text{reg}}/G$ is a genuine Riemannian manifold, so that one can as usual employ the associated Laplacian in defining the Hamiltonian on the relevant Hilbert space, $L^2(Q_{\text{reg}}/G)$.

[68]See Baker and Mulay (1995) and Landsman (1998b). It is also perhaps worth noting that a reduced phase space, considered as a singular topological space, does not automatically come equipped with a notion of smoothness—although it does inherit one if it arises via reduction from a manifold; see Pflaum (2001a, §4.4; 2001b, §5).

[69]See Tanimura and Iwai (2000) for an approach to the quantization of the singular reduced $n$-body problem. The quantization of singular quotients is an active area of mathematical research. See, e.g., Huebschmann (2002), Pflaum (2002), and the papers in Landsman *et al.* (2001).

[70]Jackiw observes (1984, 257) that the unconstrained theory admits a straightforward quantization at the heuristic level.

But if we *do* choose to impose these constraints, the dynamics of the resulting constrained theories will very likely receive different readings: the Yang-Mills case is almost always viewed intrinsically while the particle case is usually given an extrinsic construal. Why do the standard treatments of these constraints differ?

*A solution.* At this point it is helpful to consider the complementary Lagrangian formalism—and the Legendre transform relates the Lagrangian formulation of a theory to the corresponding Hamiltonian formulation. Roughly speaking, a *regular* Lagrangian leads to equations of motion with a well-posed initial value problem and to a *Legendre transform* that leads to an ordinary Hamiltonian system, while a *degenerate* Lagrangian leads to equations of motion with an ill-posed initial value problem and to a *Legendre transform* that leads to a constrained Hamiltonian system (whose equations of motion may also have an ill-posed initial value problem).[71]

We are faced with two approaches to Lagrangian mechanics (for regular Lagrangians, these lead to equivalent equations of motion for systems in the intersection of their domains of applicability). In the first, one considers a Lagrangian, $L$, defined on a tangent bundle, $TQ$ (Marsden and Ratiu (1994, §§7.1–3 and 7.7)). $L$ determines a Legendre transform $FL : TQ \rightarrow T^*Q$, which can be used to pull back the canonical symplectic form from $T^*Q$ to $TQ$, allowing one to solve for the vector field on $TQ$ determined by the Lagrangian energy. This recipe leads to a non-degenerate form on $TQ$ and well-behaved dynamics iff the Lagrangian $L$ is regular. And the Lagrangian is always regular when—as in our cases—it is just the difference between a kinetic energy determined by a (weak) Riemannian metric on $Q$ and a potential energy defined upon $Q$. So this framework does not serve to ground a distinction between the status of the constraints in the particle case in the Yang-Mills case.

A second sort of Lagrangian approach views fields as sections of bundles over spacetime, and determines the dynamics via the Euler-Lagrange equations (histories of finite dimensional systems are sections of bundles over time). If one adopts this approach to Yang-Mills theories, and works not with the Lagrangian of the previous paragraph, but with the traditional Yang-Mills Lagrangian—given by integrating $\langle\langle \mathcal{F} \rangle\rangle$ over spacetime, where $\mathcal{F}$ is the curvature associated

---

[71]The precise meaning of the italicized terms varies from context to context. For details, see Henneaux and Teitelboim (1992 Chapters 1 and 2), Kosmann-Scharwbach (1985), and Marsden and Ratiu (1994, Chapter 7). Earman (2002) provides a synthetic overview of this complicated territory from one perspective.

with the four-potential—one finds under-determined Euler-Lagrange equations and a Legendre transform that gives rise to the Hamiltonian constraint surface for our Yang-Mills theory.[72] So it becomes entirely plausible that our two constrained Hamiltonian systems deserve to be treated differently if this second Lagrangian formulation of Yang-Mills theories is taken as fundamental.

*Another solution.* The extrinsic approach to the dynamics on the constraint surface relies essentially upon the behavior of the Hamiltonian and the symplectic structure off of the constraint surface—and hence will be suspect whenever there is reason to question whether states off of the constraint surface represent genuine physical possibilities. In this situation the intrinsic dynamics will naturally be preferred. This observation provides some insight into our puzzle.

For in the particle mechanics case, orthodoxy endorses the possibility of rotating closed systems, and the correctness of the usual Newtonian account of their dynamics. This explains, as well as do considerations involving the Legendre transform, why the orthodox majority accepts an extrinsic approach to the Newtonian constraint surface. But this new consideration offers us something more: a clue as to the source in classical relationalism about space and motion—with its traditional wariness of a rotating universe—of the intuitions of the those who prefer an intrinsic reading of the dynamics.[73]

In the Yang-Mills case, meanwhile, there is a much less contentious reason to deny that states violating the constraint represent genuine physical possibilities. This may sound implausible—after all, in some contexts one *does* assign $D_A E$ non-zero values. Indeed, in one familiar form, Maxwell's equations include div $E = \rho$, for charge density $\rho$; an analogous procedure is followed in the non-abelian case. Now, this amounts to working off of the constraint surface in order to study the field dynamics in the presence of external sources painted onto spacetime independently of the behavior of the fields. This is, of course,

---

[72]See, e.g., Arms (1980) for this formulation and the corresponding Legendre transform. See Kosmann-Schwarzbach (1985) for the relations between the various Lagrangian and Hamiltonian formalisms. See Kosmann-Scharzbach (1987) or Olver (1993) for precise statements of Noether's second theorem, which, roughly speaking, says that a Lagrangian field theory has under-determined Euler-Lagrange equations iff it possesses an infinite dimensional family of variational symmetries parameterized by arbitrary functions on spacetime.

[73]One need not think of the two proposed solutions as being in competition: relationalists will be attracted to Lagrangian formulations of particle mechanics along the lines of that of Lynden-Bell (1995), where the constraints are enforced by Lagrange multipliers—resulting in a singular Lagrangian, under-determined Euler-Lagrange equations, and a Legendre transform which maps the velocity phase space onto the Hamiltonian constraint surface.

an *ad hoc* maneuver—if one wants to study Yang-Mills with sources *honestly*, one must introduce matter which not only acts upon the field but is also acted upon by it. And when one pursues this upright course, one ends up with a constraint which is a direct analog of the usual Gauss constraint—the null directions of the constraint surface correspond to the infinitesimal generators of gauge transformations.[74] Under this more fundamental approach, there is no physical interpretation for points lying off of the constraint surface—and so we have an excellent reason to prefer an intrinsic reading of the theory.

## 12. Ontology?

It is easy to imagine people who find holonomies and Wilson loops repugnant, and express their dismay by asking whether it is being seriously proposed that *spacetime non-local quantities* describe physical reality. But is hard to know what to do with the worry in this form—after all, holonomies and Wilson loops are well-defined quantities on the spaces of states of the standard formulations of Yang Mills theories. If it is accepted that these theories describe reality, doesn't it follow that the quantities in question are as real as any others?

Consider some arbitrary quantity associated with a classical system of particles— say, $\pi$ raised to the power of the square-root of the magnitude of total angular momentum. Imagine someone remarking that while this is a well-defined function on the phase space of the system, no one ought to accept such a monster as *real*. This seems like mere mistake: if it is conceded that the components of the system have positions and momenta, then it follows that the system possesses some determinate property from the determinable picked out by our function. If there is a defensible strategy in the neighborhood, it is to grant that this quantity is indeed *real*, but to insist that it is *derivative*.

Let, analogously, us ask whether holonomies enjoy a primary or a derivative status in the proper understanding of Yang-Mills theories. What is the alternative to taking holonomies as basic? Recall how we were led to them in the first place: they determine, up to a gauge transformation, the geometry of a connection on a principal $G$-bundle $P \to M$, together with the topology of the bundle. Let us for short simply say that the holonomies *determine the geometry*

---

[74]See Marsden and Weinstein (1982) for the coupling of the Maxwell field to a plasma, Weinstein (1982) for the coupling of the Maxwell field to a charged quantum particle; Śniatycki *et al.* (1996) for the coupling of the Yang-Mills field to the Dirac field; and Kuperschmidt (1992) for a large number of related constructions.

*of a bundle.* Now it seems we could insist that what Yang-Mills theories are fundamentally about is the geometry of bundles—a solution to the equations of the theory describes how this geometry changes in time. Stipulating, for convenience that spacetime is the sort of thing which carries a Lorentzian metric, we would then want to say that the theorydescribes the geometry of a world dim $G$ dimensions larger than spacetime.[75] It would, of course, remain true that such a geometry is determined by the values of certain functions on the space of loops in spacetime.

Now, it seems, we have a genuine dispute. For we can easily imagine philosophers who take themselves to live in a world correctly described by some non-abelian Yang-Mills theory, and for whom the question whether holonomies ought to be considered primary or derivative assumes the form of the question whether the world should be thought of as exhausted by spacetime, or as being some dimensions larger, with a geometry described by a Yang-Mills theory. Let us call our two contending interpretations the *holonomy interpretation* and the *connection interpretation*, respectively.

We can imagine how the dispute will go.

- The partisans of the connection interpretation will hope for a quick victory under the banner of Quineanism. They will maintain that, although there is a nice mathematical correspondence between the constrained and reduced formulations of any Yang-Mills theory—and hence a sense in which, to a first approximation, the two formulations must be compatible with the same (deterministic) interpretations of the theory—, they can nonetheless make out a finer sense in which the connection interpretation sits more nicely with the constrained formulation than does the holonomy interpretation. This motivates the suggestion that we ought to make like "Quine" and settle the ontological dispute by seeing which formulation is more useful scientifically—declaring victory for the connection interpretation once it is observed that holonomies are unwieldy and physicists are quite wedded to the variables of the constrained approach.

- The other party can afford to grant everything up to the appearance of "Quine"; then remark that they want nothing to do with this character

---

[75]It is important to distinguish the strategy here considered from the quite distinct Kaluza-Klein strategy. The latter involves *spacetimes* (in the sense stipulated in the text) of dimension greater than four, and (in field theory, though perhaps not in string theory) is beset by technical difficulties at the quantum level.

if he seeks to extinguish their right to draw a distinction between a computationally convenient set of variables and a perspicuous formulation of a theory. By way of countering the impression that holonomies are too intractable to play any real role in physics, they may also want to add that these variables in fact play an important and increasing role in the quantum theory.[76]

- At this stage, the debate will likely come down to metaphysical intuitions. The holonomy crowd may mount a counterattack, pointing out that their interpretation is the more parsimonious. And if even the connection interpretation cannot not offer locality in *spacetime*, why baulk at the non-local nature of holonomies? The other side can rejoin that it is a shortcoming of holonomies that, satisfying non-trivial constraints, they violate certain combinatorial principles (if *this* is possible and *that* is possible, then so is their juxtaposition).[77] And away we go. . .

We can imagine all of this—with relish even. But there is, I think, a serious question about whether we ought to let this dispute spill over into *our* world.

My worry here is not simply that our world is not fundamentally described by classical Yang-Mills theories. Indeed, I think we are in general entirely reasonable in debating the correct interpretation of less-than-fundamental theories. Understanding how this could be so is an outstanding philosophical problem. But this much, I suppose, is uncontentious: judgments about the interest and correctness of interpretations of theories which are (in the strictest sense) false must rest ultimately upon judgments about the extent to which various interpretations of a given theory contribute to, and integrate smoothly with, our understanding of the world. Here the following sorts of considerations play a role: background metaphysical commitments and hopes; judgments about the relative perspicuity of various alternative formulations of the theory that we are interested in, and about the links between variant formulations and competing interpretations; and considerations—operating at the technical, conceptual, and metaphysical levels—that arise when we consider how our theory is related to neighboring theories, both more and less fundamental.

---

[76]In addition to their starring role in the loop quantum gravity program, holonomies are a crucial ingredient in standard definitions of quark confinement, and figure in such highbrow topics as holography (on this last point, see Susskind and Toumbas (2000) and Rehren (2000)).

[77]This suggestion is due to Frank Arntzenius.

Now, classical non-abelian Yang-Mills theories are very unusual— *perhaps unique?*—among our menagerie of physical theories in making virtually no *direct* contribution to our understanding of the world. They are subservient to their quantum counterparts: their formulation under-girds the construction of quantum gauge theories; the study of their solutions provides a toe-hold for the construction of solutions of the quantum equations; it is always the quantum theories, never the classical ones, that play a role in applications and predictions. Classical Yang-Mills theories contribute to our knowledge of the world only via the contributions of their deeper, quantum cousins. At the same time, the vast bulk of our understanding of quantum gauge theories derives from perturbative calculations—which rely upon gauge fixing schemes and hence cast little light on the interpretative questions surrounding the gauge invariance of quantum and classical gauge theories.[78]

I conclude that clarification of our interpretative questions likely awaits deeper investigations of non-perturbative aspects of quantum gauge theories— for therein lies our best hope of anchoring interpretative claims about the classical theories. It is, however, far from clear that such investigations would ever give any traction to our present relatively simple-minded questions about the proper understanding of classical gauge freedom.

One final point: the same features of the role of classical Yang-Mills theories within our overall physics that generate the mushiness surrounding interpretative questions generate a similar mushiness surrounding questions of formulation—and this latter phenomenon is potentially unhealthy for the connection interpretation. Prior to quantizing a field theory, it is normally necessary to complete the configuration space—replacing the space, $Q$, of smooth solutions of the classical equations by some space, $\bar{Q}$, of distributional solutions; informally speaking, the states of the resulting quantum theory are wave functions over $\bar{Q}$ rather than $Q$.[79] In our case, there is a straightforward and widely investigated means of effecting this strategy: replace the configuration space, $\mathcal{A}/\mathcal{G}_*$, which consists of *smooth* assignments of holonomies to loops, by a space in which distributional assignments are allowed. In making this transition, we

---

[78] *Gauge fixing*: the goal is to circumvent gauge freedom by (smoothly) choosing a single representative of each gauge equivalence class on the constraint surface. If successful, the space of representatives is a symplectic manifold isomorphic to the reduced phase space—indeed, gauge fixing amounts to selecting a section of the principal $\mathcal{G}_*$ bundle $\Gamma \to T^*(\mathcal{A}/\mathcal{G}_*)$. This is always possible locally, but is often globally impossible. Indeed, a global obstruction can arise even when $P \to M$ is trivial and $M$ simply connected; see Singer (1978).

[79] See Wald (1994, §3.2) for this sort of construction in the linear case.

lose our interpretation of the configuration space as parameterizing the geometries of a $(4 + \dim G)$-dimensional space (see Lewandowski (1993)). Now, given the peculiar role that classical gauge theories play, we have scant grounds for insisting that one or another candidate formulation of the classical theory is the *correct* one; it becomes especially difficult to insist that the classical theory is a theory of smooth rather than distributional solutions; and the connection interpretation is in peril.

## 13. What is the Point of Gauge Freedom?

Most classical physical theories have standard formulations as simple mechanical systems: their dynamical variables split into (generalized) position variables and (generalized) momentum variables; their dynamics are given by Hamilton's equations for a Hamiltonian which is the sum of a kinetic term and a potential term.

A few important classical theories have standard formulations as gauge theories (in the sense delineated in §§2–4), in which the dynamical variables are separable into position variables and momentum variables and the Hamiltonian assumes the standard form, but the specification of initial data is subject to certain constraints and the initial value problem is well-posed only up to a time-dependent group transformation. In this situation, we can always choose local coordinates $(q^1, \ldots; p_1, \ldots; r^1, \ldots)$ so that the evolution of the $r^i$ is arbitrary, while the evolution of the $q^i$ and $p_i$ are given by a set of Hamilton's equations.

Why do some, but only some, theories assume this strange form?

One of the important lessons of the discussion of §§2–5 above is that every such gauge theory can be reduced, yielding a simple mechanical system as the reduced theory.[80] Obversely, every (well-behaved) simple mechanical system can be *enlarged* to yield a gauge theory—that is, for any given simple mechanical system, it is (almost always) possible to construct a gauge theory whose reduced theory is the given system.[81]

---

[80]There also exist physically interesting constrained theories—such as general relativity in its standard 3+1 Hamiltonian formulation—that lie outside of our class but admit a variety of reduction.

[81]This follows for finite dimensional systems from the results of Gotay and Tuynman (1988, 1991); enlargements constructed in this manner may well be rather dull. But a given simple mechanical system can admit a number of distinct enlargements; see Guillemin and Sternberg (1990, §§8–12) for interesting enlargements of the Kepler problem (i.e., the physics of a planet

These observations show that each of our theories normally formulated as a simple mechanical system *could* be reformulated as a gauge theory, while each of our theories normally formulated as a gauge theory *could* be reformulated as a simple mechanical system. So our question becomes: Why do we tend to prefer to formulate some of our theories as gauge theories, and others as simple mechanical systems?

For a given theory, there are two virtues that we might appeal to in explaining our preference for one sort of formulation over the other: *convenience* and *perspicuity*.

I believe that we must appeal to convenience alone, so long as we restrict ourselves to considerations drawn from the classical domain. For to say that the constrained formulation of a classical Yang-Mills theory is more perspicuous than the corresponding reduced formulation is to say that the former gives us deeper insight into what the classical theory is telling us about the world. It was the burden of the previous section to show that, in the present state of play, purely classical considerations are incapable of underwriting any such judgement concerning Yang-Mills theories.

It is not difficult to discern senses in which the constrained formulations of the theories discussed in §§6–8 are more convenient than their reductions.

- In particle mechanics, we start with a configuration space which is just (a subset of) $3N$ dimensional Euclidean space. In Yang-Mills theory we start with the space of connections on some principal bundle over physical space. In both cases, we have a linear (or, at worst, affine) space equipped with a flat Riemannian metric. And in both cases, the corresponding reduced configuration space is a non-linear space carrying a curved Riemannian geometry.[82] There is, then, a considerable increase in complexity in moving from the standard configuration space to the reduced configuration space— an increase which, as we saw in §10, is only increased if we allow singular reduced spaces.

- In these cases one has quite straightforward parameterization of the extended configuration spaces (positions relative to an inertial frame, values of the vector potential) while the available parameterization of the reduced configuration space are quite awkward (both relative distances and

orbiting a fixed sun).

[82]See Littlejohn and Reinsch (1997, §IV.C) and Babelon and Viallet (1981), respectively.

holonomies are grossly overcomplete in physically interesting cases; the nonlocality of the latter pose serious technical difficulties).

Thus in our examples there is a considerable gain in mathematical tractability in working with the extended configuration space rather than the reduced configuration space.[83] It is no surprise that these theories were first written down as constrained theories with gauge freedom—nor that it continues, for the most part, to be worthwhile to put up with this gauge freedom rather than struggle with the conceptually simpler, but technically unpleasant, reduced formulations of the theories.

Now, one hopes that there lurks, somewhere, a deeper account of the grounds of this sort of gain in convenience—relative to which the considerations mentioned just now will appear quite superficial. All that I want to insist upon here is that pragmatic considerations—mathematical elegance, tractability, and convenience—will drive any such explanation of the importance of gauge freedom in Yang-Mills theories that functions exclusively at classical level.

But attention to quantum considerations could easily change this situation. For instance, our discussion of §§9 and 10 suggests the following possibilities: a constrained theory may admit quantizations which do not arise as quantizations of the corresponding reduced theory; the quantization of a singular reduced theory may require exogenous structure which is most naturally viewed as deriving from a classical or quantum constrained theory.[84] If one of these scenarios were to occur in a physically important case, it would be very natural to conclude that the constrained formulation of the classical theory enjoyed an advantage in perspicuity over the reduced formulation, in that it offered us deeper insight into what the classical theory told us about our world. And if the constrained formulation (in either intrinsic or extrinsic form) enjoyed a closer association with one or another ontological picture than did the reduced formulation, then these quantum considerations would appear to have considerable interpretative consequences at the classical level.[85] In this case, our analysis of the interest

---

[83]See Kazhdan *et al.* (1978) and Marsden and Weinstein (1983) for other important sorts of example in which simplification is achieved through the introduction of additional variables.

[84]Passing to the reduced phase space may result in loss of information encoded in the constrained formulation. It might turn out to be desirable to work with a subtler reduced object that retains this information—perhaps by taking the reduced space to be a Lie groupoid rather than a mere manifold (see Weinstein (1996) for an introduction to groupoids and their uses).

[85]In roughly the same sense that the quantum Aharonov-Bohm effect gives us an insight into the ontological import of the classical vector potential; see Belot (1998).

of gauge freedom would carry us far beyond the merely pragmatic domain of mathematical tractability.[86]

**Acknowledgements**

# References

[1] Abraham, R. and J. Marsden (1985). *Foundations of Mechanics.* Cambridge, MA: Perseus.

[2] Abraham, R., J. Marsden, and T. Ratiu (1988). *Manifolds, Tensor Analysis, and Applications.* New York: Springer-Verlag.

[3] Alekseevsky, D., A. Kriegl, M. Losik, and P. Michor (2001). The Riemannian Geometry of Orbit Spaces. The Metric, Geodesics, and Integrable Systems. LANL pre-print math.DG/0102159.

[4] Arms, J. (1981). The Structure of the Solution Set for the Yang-Mills Equations. *Mathematical Proceedings of the Cambridge Philosophical Society, 90,* 361–372.

[5] Arnold, V.I. (1989). *Mathematical Methods of Classical Mechanics.* New York: Springer-Verlag.

[6] Ashtekar, A., L. Bombelli, and O. Reula (1991). The Covariant Phase Space of Asymptotically Flat Gravitational Fields. In M. Francaviglia (ed.), *Mechanics, Analysis, and Geometry: 200 Years after Lagrange* (pp. 417–50). Amsterdam: Elsevier.

[7] Ashtekar, A. and G. Horowitz (1986). On the Canonical Approach to Quantum Gravity. *Physical Review* D, 26, 3342–3353.

---

[86]Of course, there is no reason to think that such considerations, were they to obtain, would vindicate our current judgements about which theories ought to be formulated as gauge theories.

[8] Asorey, M. (1999). Maximal Non-Abelian Gauges and Topology of Gauge Orbit Space. *Nuclear Physics* B, 551, 399–424.

[9] Babelon, O. and C. Viallet (1981). The Riemannian Geometry of the Configuration Space of Gauge Theories. *Communications in Mathematical Physics,* 81, 515–525.

[10] Baker, G. and S. Mulay (1995). Geometrical and Analytical Aspects of Anyons. *International Journal of Theoretical Physics,* 34, 2435–2451.

[11] Barbour, J. and B. Bertotti (1982). Mach's Principle and the Structure of Dynamical Theories. *Proceedings of the Royal Society of London* A, 382, 295–306.

[12] Barrett, J. (1991). Holonomy and Path Structures in General Relativity and Yang-Mills Theory. *International Journal of Theoretical Physics,* 30, 1171–1215.

[13] Bates, L. and E. Lerman (1997). Proper Group Actions and Symplectic Stratified Spaces. *Pacific Journal of Mathematics,* 181, 201–229.

[14] Bates, S. and A. Weinstein (1997). *Lectures on the Geometry of Quantization.* Providence, RI: American Mathematical Society.

[15] Batterman, R. (2002). Falling Cats, Parallel Parking, and Polarized Light. PITT-PHIL-SCI pre-print 00000583.

[16] Belot, G. (1998). Understanding Electromagnetism. *British Journal for the Philosophy of Science,* 49, 531–555.

[17] Crnkovic, C. and E. Witten (1987). Covariant Description of Canonical Formalism in Geometrical Theories. In S. Hawking and W. Israel (eds.), *Three Hundred Years of Gravitation*(pp. 676–84). New York: Cambridge University Press.

[18] Cushman, R. and L. Bates (1997). *Global Aspects of Integrable Systems.* Boston: Birkhaüser.

[19] Diacu, F. and P. Holmes (1996). *Celestial Encounters.* Princeton: Princeton University Press.

[20] Dirac, P. (1964), *Lectures on Quantum Mechanics.* New York: Yeshiva University Press.

[21] Duval, C., J. Elhadad, M. Gotay, and G. Tuynman (1990). Nonunimodularity and the Quantization of the Pseudo-Rigid Body. In J. Harnad and J. Marsden (eds.), *Hamiltonian Systems, Transformation Groups, and Spectral Transform Methods* (pp. 149–60). Montréal: Les Publications CRM.

[22] Earman, J. (2002).Gauge Matters. *Philosophy of Science,* 69: S209–S220.

[23] Emmrich, C. and H. Römer (1990). Orbifolds as Configuration Spaces of Systems with Gauge Symmetries. *Communications in Mathematical Physics,* 129, 69–94.

[24] Göckeler, M. and T. Schücker (1987). *Differential Geometry, Gauge Theories, and Gravity.* New York: Cambridge University Press.

[25] Gotay, M. (1986a). Constraints, Reduction, and Quantization. *Journal of Mathematical Physics,* 27, 2051–2066.

[26] Gotay, M. (1986b). Negative Energy States in Quantum Gravity? *Classical and Quantum Gravity,* 3, 487–491.

[27] Gotay, M. (1989). Reduction of Homogeneous Yang-Mills Fields. *Journal of Geometry and Physics,* 6, 349–365.

[28] Gotay, M., J. Isenberg, and J. Marsden (1998). Momentum Maps and Classical Relativistic Fields. Part I: Covariant Field Theory. LANL pre-print physics/9801019.

[29] Gotay, M. and G. Tuynman (1989). $\mathbf{R}^{2n}$ is a Universal Symplectic Manifold for Reduction. *Letters in Mathematical Physics,* 18, 55–59.

[30] Gotay, M. and G. Tuynman (1991). A Symplectic Analogue of the Mostow-Palais Theorem. In P. Dazord and A. Weinstein (eds.), *Symplectic Geometry, Groupoids, and Integrable Systems* (pp. 173–83). New York: Springer-Verlag.

[31] Guillemin, V. and S. Sternberg (1990). *Variations on a Theme by Kepler.* Providence, RI: American Mathematical Society.

[32] Healey, R. (2001). On the Reality of Gauge Potentials. *Philosophy of Science*, 68, 432–455.

[33] Henneaux, M. and C. Teitelboim (1992). *Quantization of Gauge Systems*. Princeton: Princeton University Press.

[34] Huebschmann, J. (1996). The Singularities of Yang-Mills Connections for Bundles on a Surface—II. The Stratification. *Mathematische Zeitschrift*, 221, 83–92.

[35] Huebschmann, J. (2002). Kähler Quantization and Reduction. LANL preprint math.SG/0207166.

[36] Isenberg, J. and J. Marsden (1982). A Slice Theorem for the Space of Solutions of Einstein's Equations. *Physics Reports*, 89, 179–222.

[37] Isham, C. (1984). Topological and Global Aspects of Quantum Theory. In B. DeWitt and R. Stora (eds.), *Relativity, Groups and Topology* II (pp. 1059–1290). Amsterdam: Elsevier.

[38] Jackiw, R. (1984). Topological Investigations of Quantized Gauge Theories. In B. DeWitt and R. Stora (eds.), *Relativity, Groups and Topology* II (pp. 221–331). Amsterdam: Elsevier.

[39] Kazhdan, D., B. Kostant, and S. Sternberg (1978). Hamiltonian Group Actions and Dynamical Systems of Calogero Type. *Communications on Pure and Applied Mathematics*, 31, 481–507.

[40] Koiller, J., K. Ehlers, and R. Montgomery (1996). Problems and Progress in Microswimming. *Journal of Nonlinear Science*, 6, 507–41.

[41] Kondracki, W. and J. Rogulski (1986). On the Stratification of the Orbit Space for the Action of Automorphisms on Connections. *Dissertationes Mathematicae*, CCL, 1–62.

[42] Kosmann-Schwarzbach, Y. (1985). On the Momentum Mapping in Field Theory. In H.-D. Doebner and J. Hennig (eds.), *Differential Geometric Methods in Physics.* (pp. 25–71). New York: Springer-Verlag.

[43] Kosmann-Schwarzbach, Y. (1987). Sur les Théorèmes de Noether. In Y. Choquet-Bruhat, B. Coll, R. Kerner, and A. Lichnerowicz (eds.), *Géométrie et Physique* (pp. 149–160). Paris: Hermann.

[44] Kuchař, K. (1988). Canonical Quantization of Generally Covariant Systems. In B. Iyer, A. Khembhavi, J. Narlikar, and C. Vishveshvara (eds.), *Highlights in Gravitation and Cosmology* (pp. 93–120). New York: Cambridge University Press.

[45] Kuperschmidt, B. (1992). *The Variational Principles of Dynamics*. Singapore: World Scientific.

[46] Landsman, N. (1998a). *Mathematical Topics Between Classical and Quantum Mechanics*. New York: Springer-Verlag.

[47] Landsman, N. (1998b). Quantization of Singular Systems and Incomplete Motions. In M. Rainer and H.-J. Schmidt (eds.), *Current Topics in Mathematical Cosmology* (pp. 256–263). Singapore: World Scientific.

[48] Landsman, N., M. Pflaum, and M. Schlichenmaier (eds.) (2001). *Quantization of Singular Symplectic Quotients*. Boston: Birkhäuser.

[49] Lang, S. (1999). *Fundamentals of Differential Geometry*. New York: Springer-Verlag.

[50] Lewandowski, J. (1993). Group of Loops, Holonomy Maps and Path Connection. *Classical and Quantum Gravity, 10*, 879–904.

[51] Littlejohn, R. and M. Reinsch (1997). Gauge Fields in the Separation of Rotations and Internal Motions in the $n$-Body Problem. *Reviews of Modern Physics, 69*, 213–275.

[52] Liu, C. (2001). Gauge Gravity and the Unification of Natural Forces. PITT-PHIL-SCI pre-print 00000364.

[53] Loll, R. (1992). Canonical and BRST-Quantization of Constrained Systems. In M. Gotay, J. Marsden, and V. Moncrief (eds.), *Mathematical Aspects of Classical Field Theory* (pp. 503–530). Providence, RI: American Mathematical Society.

[54] Loll, R. (1994). Gauge Theory and Gravity in the Loop Formulation. In J. Ehlers and H. Friedrich (eds.), *Canonical Gravity* (pp. 254–88). New York: Springer-Verlag.

45

[55] Lynden-Bell, D. (1995). A Relative Newtonian Mechanics. In J. Barbour and H. Pfister (eds.), *Mach's Principle* (pp. 172–178). Boston: Birkhäuser.

[56] Marathe, K. and G. Martucci (1992). *The Mathematical Foundations of Gauge Theories.* Amsterdam: North-Holland.

[57] Marsden, J. (1992). *Lectures on Mechanics.* New York: Cambridge University Press.

[58] Marsden, J. and R. Abraham (1970). Hamiltonian Mechanics on Lie Groups and Hydrodynamics. In S.-S. Chern and S. Smale (eds.), *Global Analysis* (pp. 237–244). Providence, RI: American Mathematical Society.

[59] Marsden, J. and T. Ratiu (1994). *Introduction to Mechanics and Symmetry.* New York: Springer-Verlag.

[60] Marsden, J. and A. Weinstein (1982). The Hamiltonian Structure of the Maxwell-Vlasov Equations. *Physica* D, 4, 394–406.

[61] Marsden, J. and A. Weinstein (1983). Coadjoint Orbits, Vortices and Clebsch Variables for Incompressible Fluids. *Physica* D, 7, 305–323.

[62] Mitter, P. (1980). Geometry of the Space of Gauge Orbits and the Yang-Mills Dynamical System. In G. 't Hooft, C. Itzykson, A. Jaffe, H. Lehmann, P. Mitter, I. Singer, and R. Stora (eds.), *Recent Developments in Gauge Theories* (pp. 265–292). New York: Plenum Press.

[63] Moncrief, V. (1980). Reduction of the Yang-Mills Equations. In P. Garcia, A. Pérez-Rendón, and J.-M. Souriau (eds.), *Differential Geometrical Methods in Mathematical Physics* (pp.276–291). New York: Springer-Verlag.

[64] Obukhov, Y. (2000). On Physical Foundations and Observational Effects of Cosmic Rotation. M. Scherfner, T. Chrobok and M. Shefaat (eds.), *Colloquium on Cosmic Rotation* (pp. 23–96). Berlin: Wissenschaft und Technik Verlag.

[65] Olver, P. (1993). *Applications of Lie Groups to Differential Equations.* New York: Springer-Verlag.

[66] Ortega, J.-P. and T. Ratiu (1998). Singular Reduction of Poisson Manifolds. *Letters in Mathematical Physics,* 46, 359–372.

[67] Pflaum, M. (2001a). *Analytic and Geometric Study of Stratified Spaces*. New York: Springer-Verlag.

[68] Pflaum, M. (2001b). Smooth Structures on Stratified Spaces. In N. Landsman, M. Pflaum, and M. Schlichenmaier (eds.), *Quantization of Singular Symplectic Quotients* (pp. 231–258). Boston: Birkhäuser.

[69] Pflaum, M. (2002). On the Deformation Quantization of Symplectic Orbispaces. LANL pre-print math-ph/0208020.

[70] Pooley, O. and H. Brown (2002). Relationalism Rehabilitated? I: Classical Mechanics. *British Journal for the Philosophy of Science*, 53: 183–204.

[71] Rajeev, S. and L. Rossi (1995). Some Rigorous Results for Yang-Mills Theories on a Cylinder. *Journal of Mathematical Physics*, 36, 3308–3319.

[72] Redhead, M. (2002). The Interpretation of Gauge Symmetry. To appear in M. Kuhlmann, H. Lyre, and A. Wayne (eds.), *Ontological Aspects of Quantum Field Theory*.

[73] Rehren, K.-H. (2000). Algebraic Holography. *Annales Henri Poincaré*, 1, 607–623.

[74] Römer, H. (1988). Singular Points in Level Sets of the Momentum Map and Quantum Theory. In K. Bleuler and M. Werner (eds.), *Differential Geometric Methods in Theoretical Physics* (pp. 307–315). Dordrecht: Kluwer.

[75] Rudolph, G., M. Schmidt, and I. Volobuev (2002), On the Gauge Orbit Space Stratification (A Review). *Journal of Physics* A, 35: R1–R50..

[76] Schmid, R. (1987). *Infinite Dimensional Hamiltonian Systems*. Naples: Bibliopolis.

[77] Singer, I. (1978). Some Remarks on the Gribov Ambiguity. *Communications in Mathematical Physics*, 60, 7–12.

[78] Singer, I. (1982). On Yang-Mills Fields. In A. Bishop, D. Campbell, and B. Nikolaenko (eds.), *Nonlinear Problems* (pp. 35–50). Amsterdam: North-Holland.

[79] Sjamaar, R. (1996). Symplectic Reduction and the Riemann-Roch Formula for Multiplicities. *Bulletin of the American Mathematical Society,* 33, 327–338.

[80] Śniatycki, J. (1999). Regularity of Constraints and Reduction in the Minkowski Space Yang-Mills-Dirac Theory. *Annales de l'Institut Henri Poincaré* Physique Théorique, 70, 277–293.

[81] Śniatycki, J. (2000). Quantization of Yang-Mills Fields Commutes with Reduction. H.-D. Doebner, J.-D. Hennig, W. Lücke, and V. Dobrev (eds.), *Quantum Theory and Symmetry* (pp. 68–75). Singapore: World Scientific.

[82] Śniatycki, J., G. Schwarz, and L. Bates (1996). Yang-Mills and Dirac Fields in a Bag, Constraints and Reduction. *Communications in Mathematical Physics,* 176, 95–115.

[83] Susskind, L. and N. Toumbas (2000). Wilson Loops as Precursors. *Physical Review* D, 61, 044001.

[84] Tanimura, S. and T. Iwai (2000). Reduction of Quantum Systems on Riemannian Manifolds with Symmetry and Application to Molecular Mechanics. *Journal of Mathematical Physics,* 41, 1814–1842.

[85] Tuynman, G. (1992). What Are the Rules of the Game Called BRST? In M. Gotay, J. Marsden, and V. Moncrief (eds.), *Mathematical Aspects of Classical Field Theory* (pp. 625–633). Providence, RI: American Mathematical Society.

[86] Vilela Mendes, R. (2000). Gauge Strata and Particle Generations. LANL pre-print hep-th/0009027.

[87] Wald, R. (1994). *Quantum Field Theory in Curved Spacetime and Blackhole Thermodynamics.* Chicago: University of Chicago Press.

[88] Weinstein, A. (1982). Gauge Groups and Poisson Brackets for Interacting Particles and Fields. M. Tabor and Y. Treve (eds.), *Mathematical Methods in Hydrodynamics and Integrability in Dynamical Systems* (pp.1–11). New York: American Institute of Physics.

[89] Weinstein, A. (1996). Groupoids: Unifying Internal and External Symmetry. *Notices of the American Mathematical Society,* 43, 744–752.

[90] Woodhouse, N. (1980). *Geometric Quantization.* New York: Oxford University Press.