

Introduction to the Yang-Mills quantum theory

R. Jackiw*

Centro de Investigación del I.P.N., Apartado Postal 14-740, Mexico 14, D.F.

A pedagogical discussion of the Yang-Mills quantum theory is presented. A somewhat unconventional description makes use of physical, quantum-mechanical ideas, rather than of formal, mathematical developments. The purpose is to highlight those aspects of the model which have been exposed in the last few years by semiclassical methods, but without using semiclassical approximations. This requires a careful treatment of the non-Abelian gauge symmetry present in the theory.

CONTENTS

I. Introduction	661
II. Description of the Yang-Mills Theory	662
III. Canonical Formalism for Gauge Theories	664
IV. Gauge Transformations, Topology, and the Vacuum Angle	665
V. Solving Gauss' Law Constraints	670
VI. Conclusion	672
Acknowledgments	672
References	672

I. INTRODUCTION

Non-Abelian gauge theories, of the kind developed by Yang and Mills (1954), have become the focus of widespread interest, owing to their central role in two currently popular models for fundamental physical processes: the "electroweak" quantum flavor dynamics¹ and the "strong" quantum chromodynamics.² We have arrived at Yang-Mills gauge theories from both physical, phenomenological as well as theoretical, abstract considerations.

The history of the gauge principle begins with the analysis by Weyl (1950a) of gravitational theory and electrodynamics. Motivated by a desire to extend this principle to isospin transformations; and evidently unaware of the interesting precursor Klein (1939), nor of simultaneous work by Shaw (1955), Yang and Mills produced the celebrated model. However, successful application of their ideas had to await the maturation of our understanding. On the one hand, the practical importance of vector meson particles, which became evident through the work of Sakurai (1960) and others, pointed again to Yang-Mills theories, since they involve vector fields. On the other, theorists learned how to quantize and renormalize non-Abelian gauge theories; and, most importantly, it was shown that the gauge symmetry can be spontaneously broken by the Goldstone-Higgs mechanism, so that the gauge fields can describe massive vector mesons, while still retaining their renormalizable interactions. This then led to the first successful description of physical reality in terms of $SU(2) \otimes U(1)$ Yang-Mills fields coupled to fermions and symmetry-breaking Higgs scalars (quantum flavor dynamics): the Weinberg-Salam model, unifying electro-

magnetic with weak interactions, where the gauge fields are identified with the massless photon, and with the hypothetical massive vector mesons mediating weak interactions (Weinberg, 1967; Salam, 1968).

The importance of Yang-Mills theory for strong interactions derives from the fact that this is the only dynamical model whose forces become negligible at short distances, a phenomenon called "asymptotic freedom" ('t Hooft, 1972; Gross and Wilczek, 1973; Politzer, 1973). Consequently, $SU(3)$ Yang-Mills fields coupled to quarks (quantum chromodynamics) appear to provide the only realistic framework that can accommodate the MIT/SLAC experiments on high-energy lepton-nucleon scattering. The further discovery by MIT/SLAC experimentalists of the J/ψ particles left few skeptical about the physical importance of non-Abelian gauge fields, which, however, are not identified with observed particles, but merely provide the "glue" that keeps the quarks bound inside hadrons, so strongly that they are permanently confined. (These ideas about strong interactions await definitive theoretical proof.)

Of course the Yang-Mills field equations have not been solved exactly, not even in the context of classical field theory. Our understanding was originally achieved by perturbative methods. The initial approaches to the quantum theory made much use of our knowledge of the completely solvable noninteracting limit, as well as of the well-understood, perturbatively solvable quantum electrodynamics, which was perceived as a simple paradigm. This approximate, perturbative development yielded a picture of the quantum theory which possesses many physically desirable features, but still leaves much uncertainty about how thoroughly successful is the account of natural phenomena.³

In an attempt to uncover further properties of the theory, we turned to nonperturbative, semiclassical approximation methods, which showed the model to possess a much richer physical content than had been heretofore appreciated.⁴ Thus while we still cannot say that we know completely the physical predictions of Yang-Mills theory and that they agree with observed phenomena in all aspects, we have in hand an excellent candidate for a model of fundamental physics.

The recent nonperturbative results have been largely

*Permanent address: Center for Theoretical Physics, Laboratory for Nuclear Science and Department of Physics Massachusetts Institute of Technology, Cambridge, Massachusetts 02139.

¹For a review, see Taylor (1976).

²For a review, see Marciano and Pagels (1978).

³For a review of Yang-Mills quantum theory, which however does not cover the recent results summarized here, see Faddeev and Slavnov (1980).

⁴For a review of the recent nonperturbative results, see Coleman (1977, 1979), Jackiw (1977), and Jackiw, Nohl, and Rebbi (1978).

viewed as additions, corrections, and modifications of the initially developed physical picture. Yet it is possible now, with the hindsight of the semiclassical analysis, to realize that the original development suffered from omissions and oversights which can be corrected without recourse to any approximation method. Rather, when closer attention is paid to the canonical derivation of the quantum theory, we can find much of the recently discovered behavior, not as a semiclassical feature of the quantum theory, but as a consequence of general quantum-mechanical principles when realized in the Yang-Mills model.

In this article, I present a pedagogical discussion of how to quantize Yang-Mills theory, with emphasis on the unexpected aspects, which, however, will be seen to arise quite naturally when the derivation is carried out in a definite way. Analogy will be drawn to familiar, quantum-mechanical examples exhibiting similar behavior in simple, well-understood settings. I shall not be following the operator route to the quantum theory based on the electromagnetic, Coulomb-gauge analogy. Nor shall I begin with the functional integral, which leaves properties of the states obscure. However, as will be seen, both these conventional approaches can be found at the end of the present development.

Few new results, save the unified framework, are offered. Nevertheless, I trust that an audience including not only particle physicists, but also relativists and mathematicians, will find the presentation instructive. I suspect that a similar development can be carried through for the recently posited CP^N models, as well as for gravity theory. This should be done since it may very well produce new insight into these difficult-to-solve examples. Some initial investigations already exhibit this (Deser, Duff, and Isham, 1980; Friedman and Sorkin, 1980; Isham, 1980).

II. DESCRIPTION OF THE YANG-MILLS THEORY

The basic dynamical variables of the Yang-Mills theory are the vector potentials A_μ^a , carrying space-time index μ , and internal symmetry index a .⁵ Although several Lie groups occur in physical applications [$SU(2) \times U(1)$ for the electroweak interactions, $SU(3)$ for the strong, larger ones in speculative models that attempt to unify the strong with the electroweak], we shall confine the discussion to $SU(2)$, as the phenomena which we wish to highlight depend only on the non-Abelian nature of the group. Thus a ranges over 1, 2, and 3. The component notation, A_a^μ , will be used interchangeably with the matrix notation, $A^\mu = A_a^\mu (\sigma^a / 2i)$, where $\sigma^a =$ Pauli matrices. The Yang-Mills fields $F_{\mu\nu}$ are related to the potentials A_μ by

$$F_{\mu\nu} = \partial_\mu A_\nu - \partial_\nu A_\mu + g[A_\mu, A_\nu]. \tag{1}$$

Here g is the coupling constant. The fields satisfy an equation of motion

$$\mathcal{D}_\mu F^{\mu\nu} = 0 \tag{2}$$

$$\mathcal{D}_\mu = \partial_\mu + g[A_\mu, \] \tag{3}$$

⁵The diagonal metric which we use, $g^{\mu\nu}$, has diagonal entries (1, -1, -1, -1); $\epsilon^{\mu\nu\alpha\beta}$ is the totally antisymmetric Levi-Civita tensor with $\epsilon^{0123} = 1$. Throughout \hbar and c are set to unity.

which is also the Euler-Lagrange equation for the condition that the action I

$$I = \int d^4x \mathcal{L} \tag{4a}$$

$$\mathcal{L} = \frac{1}{2} \text{tr} F^{\mu\nu} F_{\mu\nu} \tag{4b}$$

be stationary against arbitrary variations of the potentials. The dual fields $*F^{\mu\nu}$, defined by

$$*F^{\mu\nu} = \frac{1}{2} \epsilon^{\mu\nu\alpha\beta} F_{\alpha\beta}, \tag{5}$$

satisfy the Bianchi identity, which is a consequence of the definitions (1), (3), and (5)

$$\mathcal{D}_\mu *F^{\mu\nu} = 0. \tag{6}$$

The theory is gauge invariant, i.e., it is invariant under the transformation

$$A_\mu \rightarrow U^{-1} A_\mu U + g^{-1} U^{-1} \partial_\mu U, \tag{7}$$

where U is an element of the group—a 2×2 , space-time-dependent, unitary matrix with unit determinant. Indeed the equations were constructed so that this invariance would be present. In infinitesimal form the symmetry transformation, θ_a .

$$U = e^\Omega \approx I + \Omega, \Omega^\dagger = -\Omega \tag{8}$$

$$\delta A_\mu = g^{-1} \mathcal{D}_\mu \Omega \tag{9}$$

with Ω being related to the infinitesimal, local parameters of the transformation, θ_a .

$$\Omega(x) = \theta_a(x) \left(\frac{\sigma^a}{2i} \right). \tag{10}$$

The field strengths are not gauge invariant; rather they are gauge covariant. They transform homogeneously under gauge transformations, in contrast to the potentials which follow the inhomogeneous transformation law (7)

$$F^{\mu\nu} \rightarrow U^{-1} F^{\mu\nu} U, \tag{11}$$

$$\delta F^{\mu\nu} = [F^{\mu\nu}, \Omega]. \tag{12}$$

Note that the various equations which the fields satisfy cannot be written in terms only of $F^{\mu\nu}$, as is the case for the Abelian Maxwell theory, since the covariant derivatives \mathcal{D}_μ involve the potential A_μ . This reflects the fact that the non-Abelian potentials play a more fundamental role in the Yang-Mills theory than do their Abelian counterparts in the Maxwell theory. The physical content is not coded completely in the field strengths, which are not even gauge invariant. An object that does contain all the gauge-invariant information is the nonintegrable phase factor (Wu and Yang, 1975)

$$P(C) = \text{tr} P \exp \oint_C dz^\mu A_\mu(z). \tag{13}$$

The integration is path ordered by the symbol P , and it runs over a closed contour C . The resulting quantity is gauge invariant, but path dependent. Knowledge of $P(C)$ for arbitrary contours allows one to reconstruct the potentials, up to gauge transformations. At the present time there are many attempts to write the theory solely in terms of these objects, but this pro-

gram has thus far not been completed, and will not be discussed here.

It is also interesting to discuss how the potentials transform under coordinate transformations.

$$x^\mu \rightarrow \tilde{x}^\mu(x). \tag{14a}$$

A conventional statement is that the transformed potentials, evaluated at the transformed point, $\tilde{A}_\mu(\tilde{x})$, satisfy a formula appropriate to a coordinate vector quantity. This is the familiar transformation law which is postulated in Riemannian geometry.

$$\tilde{A}_\alpha(\tilde{x}) \frac{\partial \tilde{x}^\alpha}{\partial x^\mu} = A_\mu(x). \tag{14b}$$

When the coordinate change is represented infinitesimally

$$\begin{aligned} \tilde{x}^\mu &= x^\mu + \delta x^\mu \\ \delta x^\mu &= -f^\mu(x) \end{aligned} \tag{14c}$$

then Eq. (14b) reduces to a formula involving the Lie derivative L_f .

$$\begin{aligned} \tilde{A}_\mu(x) &= A_\mu(x) + \delta_f A_\mu(x) \\ \delta_f A_\mu &= f^\alpha \partial_\alpha A_\mu + (\partial_\mu f^\alpha) A_\alpha \equiv L_f A_\mu. \end{aligned} \tag{14d}$$

[Within this general expression, we recognize familiar special cases. For example, for translations f^α is a constant, a^α , and $\delta_{\text{translation}} A_\mu = a^\alpha \partial_\alpha A_\mu$, which is of course the standard result. Similarly for Lorentz transformations $f^\alpha = \omega^{\alpha\beta} x_\beta$, $\omega^{\alpha\beta} = -\omega^{\beta\alpha}$, and $\delta_{\text{Lorentz}} A_\mu = -\frac{1}{2} \omega^{\alpha\beta} [(x_\alpha \partial_\beta - x_\beta \partial_\alpha) A_\mu + g_{\alpha\mu} A_\beta - g_{\beta\mu} A_\alpha]$, which again is the usual Lorentz transformation law for a vector field.]

However, because we are dealing with a gauge theory, there is the possibility of modifying the conventional transformation law (14d), which was invented for arbitrary vector fields, not necessarily gauge fields. Observe that the Lie derivative in (14d) may also be written as

$$\delta_f A_\mu = f^\alpha F_{\alpha\mu} + \mathcal{D}_\mu (f^\alpha A_\alpha). \tag{15}$$

Since the last term in (15) is a gauge transformation, which certainly may be adjoined at will to any transformation, we may adopt as a coordinate transformation law, instead of (14d), the following gauge-covariant expression (Jackiw, 1979):

$$\delta_f A_\mu = f^\alpha F_{\alpha\mu}. \tag{16}$$

The gauge-covariant transformation law (16) of course produces no conceptual change from the conventional law (14d). Nevertheless it allows a more elegant discussion of various formal questions concerned with Noether's theorem, supersymmetry, fiber bundles, and the like. These are beyond the scope of the present review and are summarized in the literature.⁶ [The finite transformation corresponding to the infinitesimal Eq. (16) is nontrivial. It involves the nonintegrable phase factor (Jackiw, 1979).]

One verifies when f^α is a conformal Killing vector, i.e., when f^α satisfies the equation

⁶For a review, and application to the study of symmetry and invariance in a gauge theory, see Jackiw (1980).

$$\partial_\alpha f_\beta + \partial_\beta f_\alpha - \frac{1}{2} g_{\alpha\beta} \partial_\mu f^\mu = 0 \tag{17}$$

that the transformation (16) [or (14d)] is a symmetry operation for the action (4).

$$\delta_f I = 0. \tag{18}$$

The corresponding conserved current, which follows from an application of Noether's theorem, is

$$J^\mu = \theta^{\mu\nu} f_\nu \tag{19}$$

where $\theta^{\mu\nu}$ is a symmetric and traceless energy-momentum tensor.⁷

$$\theta^{\mu\nu} = 2 \text{tr} (F^{\mu\alpha} F_\alpha^\nu - \frac{1}{4} g^{\mu\nu} F^{\alpha\beta} F_{\alpha\beta}). \tag{20}$$

This is just the conformal symmetry of the model. It gives rise to the 15-parameter $O(4, 2)$ invariance group of the Yang-Mills theory, corresponding to the 15-parameter solution to (17).

$$f^\alpha = a^\alpha + \omega^{\alpha\beta} x_\beta + c x^\alpha + b^\alpha x^2 - 2x^\alpha x \cdot b. \tag{21}$$

Here a^α , c , b^α , and $\omega^{\alpha\beta}$ are constants, the latter being antisymmetric. As mentioned already, a^α describes translations, $\omega^{\alpha\beta}$ Lorentz transformations. These are the Poincaré invariances that one expects to hold in any realistic field theory. Classical Yang-Mills theory, however, possess further invariances, described by c and b^α ; these are the dilatations and special conformal transformations which are present owing to the absence of mass terms in \mathcal{L} .⁸ In the quantum theory, the dilatations and special conformal transformations are absent, owing to renormalization effects which necessarily introduce a mass scale; the symmetries acquire anomalies.⁸ But this symmetry is controllable and anomalously broken scale invariance is exploited with the help of the renormalization group. I shall not be reviewing this topic here.²

Equation (16) may also be used to discuss configurations of classical Yang-Mills fields which are invariant under a coordinate transformation. In a conventional field theory one would say that a field configuration is invariant when its infinitesimal variation vanishes. However, in a gauge theory, the concept of invariance should be extended to allow for a possible noninvariance which can be compensated by a gauge transformation. Thus according to Eqs. (9) and (16) we see that a criterion for coordinate invariance in a gauge theory is⁹

⁷When the modified transformation (16) is applied to Fermi fields, the conserved tensor generated by Noether's theorem is gauge invariant, but not symmetric. It is in fact the energy-momentum tensor which occurs in the Kibble-Sciama (Einstein-Cartan) gravitational theory with torsion; see Kibble (1961) and Sciama (1962).

⁸For a summary, see Jackiw (1972).

⁹The study of symmetry properties of gauge fields was initiated by the mathematician H. Wang; for a mathematical summary see Kobayashi and Nomizu (1963). Recent investigations by physicists include Schwarz (1977), Bergmann and Flaherty (1978), Trautman (1979), Forgács and Manton (1980), Har-nad, Shnider, and Vinet (1980), and Jackiw and Manton (1980). A discussion from the point of view of infinitesimal gauge-invariant coordinate transformations is in Jackiw (1980).

$$\delta A_\mu = f^\alpha F_{\alpha\mu} = \mathcal{D}_\mu \Phi_f \quad (22)$$

which may be interpreted as the statement that certain projections of a coordinate invariant field strength are total (covariant) derivatives of a scalar field.¹⁰

The quantities Φ_f have the following physical significance. Consider a "matter" system of particles or fields (either classical or quantum-mechanical) moving in space and (possibly) undergoing self-interactions. As a consequence of various invariances of the matter dynamics against specific coordinate transformations, various constants of motion govern the time-evolution. Examples are the energy, arising from time-translation invariance, the angular momentum coming from rotational invariance, etc. When the same system is put in a prescribed, external, classical gauge field, the constants will in general disappear, since a generic external field breaks all invariances. However, when the external field is itself symmetric, in the sense (22), the constants remain, but in modified form. Specifically, if the matter-gauge field interaction is described by the Lagrange density $-j_\mu^a A_\mu^a$, where j_μ^a is the matter current, then the constant of motion C^f , in the presence of the symmetric external potential A_μ^a which satisfies Eq. (22), is

$$C^f = C_{\text{matter}}^f - \int dr \rho_a \Phi_f^a. \quad (23)$$

Here f labels the coordinate transformation which gives rise to the symmetry in question; the first term on the right-hand side is the matter contribution to the total conserved quantity; the second is the contribution of the external gauge potential, with ρ_a being the matter "charge" density, $\rho_a = j_0^a$ (Jackiw and Manton, 1980).

A special case of (23) is a familiar result from classical magnetic monopole theory in electrodynamics. A magnetic monopole potential

$$\mathbf{A}: A^r = 0, \quad A^\theta = 0, \quad A^\phi = -(g/r) \text{ctn}\theta \quad (24a)$$

is not manifestly spherically symmetric, but gives rise to spherically symmetric physics since the magnetic field possesses that property.

$$\mathbf{B} = \nabla \times \mathbf{A} = g(\mathbf{r}/r^3). \quad (24b)$$

Formula (22) in this application to Abelian gauge fields reads

$$\begin{aligned} f^\alpha: f^0 = 0, \quad \mathbf{f} = \mathbf{r} \times \mathbf{n}; \quad \mathbf{n} = \text{axis of rotation} \\ (\mathbf{r} \times \mathbf{n}) \times \mathbf{B} = \nabla \Phi \\ \Phi = g\mathbf{n} \cdot \mathbf{r}/r. \end{aligned} \quad (25)$$

Correspondingly, according to Eq. (23), the constant of motion—the angular momentum—of a charged point

¹⁰Equation (22) shows that a scalar field arises naturally from a (symmetric) gauge field. This result and its generalization have been used in attempts to construct from gauge fields the Higgs fields that are used in quantum flavor-dynamical models. One begins with gauge fields in a space-time of dimensionality greater than the physical four, but requires invariance against some coordinate transformations in the additional dimensions. One thus arrives at a Yang-Mills theory in four dimensions supplemented by scalar fields that play the role of Higgs fields. For this modern reprise of the Kaluza-Klein idea, see Fairlie (1979), Manton (1979), and Mayer (1980).

particle $[\rho(\mathbf{r}) = e\delta(\mathbf{x} - \mathbf{r})]$ moving in this background field, has, in addition to the usual, kinematical matter contribution $\mathbf{x} \times \mathbf{p}$, a term coming from the external gauge potential,

$$\mathbf{J} = \mathbf{x} \times \mathbf{p} - g e (\mathbf{x}/x). \quad (26)$$

This is the celebrated formula of Poincaré for the total angular momentum of a charged particle in a magnetic monopole field.

For a derivation of the above and a discussion of further developments on symmetry and invariance in gauge theories, the reader is referred to the literature.^{6,9}

III. CANONICAL FORMALISM FOR GAUGE THEORIES

We now turn to a derivation of the canonical theory. This development will be deemed successful if we can identify canonical variables, coordinates and momenta, define a Hamiltonian, and regain the field equations (1) and (2) as Hamiltonian equations. Since we are dealing with a gauge-invariant theory, we expect complications—a straightforward approach will fail, as it also fails in electrodynamics. When careful attention is paid to the subtleties of gauge invariance, it will be possible to identify in the formal, canonical development many of the features which have been recently exposed by semiclassical reasoning.

We look for the canonical momentum, but immediately we encounter the familiar gauge theory problem that the momentum conjugate to A_a^0 vanishes, since \mathcal{L} does not depend on $\partial_t A_a^0$. In order to circumvent this initial obstacle to quantization, we make use of gauge invariance to set A_a^0 to zero. The canonical variables are therefore the coordinates A_a^i and their conjugate momenta

$$\Pi_a^i = \frac{\delta \mathcal{L}}{\delta \partial_t A_a^i} = F_a^{0i} = \partial_t A_a^i. \quad (27)$$

The Hamiltonian

$$H = \int dx^3 \mathcal{H} \quad (28a)$$

$$\mathcal{H} = \Pi_a^i \partial_t A_a^i - \mathcal{L}$$

coincides with the energy; see Eq. (20).

$$\begin{aligned} H &= \frac{1}{2} \int dr (\mathbf{E}_a^2 + \mathbf{B}_a^2) = \int dr \theta^{00} \\ E_a^i &= F_a^{i0} = -\Pi_a^i, \quad B_a^i = -\frac{1}{2} \epsilon^{ijk} F_{ajk}. \end{aligned} \quad (28b)$$

When the obvious canonical commutation relations (Poisson bracket relations in the classical theory) are posited

$$\begin{aligned} [A_a^i(\mathbf{r}, t), A_b^j(\mathbf{r}', t)] &= 0 \\ [E_a^i(\mathbf{r}, t), E_b^j(\mathbf{r}', t)] &= 0 \\ [E_a^i(\mathbf{r}, t), A_b^j(\mathbf{r}', t)] &= i\delta_{ab} \delta^{ij} \delta(\mathbf{r} - \mathbf{r}') \end{aligned} \quad (29)$$

one finds that the Hamiltonian equations reproduce the definition of \mathbf{E}_a .

$$\partial_t A_a = i[H, A_a] = -\mathbf{E}_a. \quad (30a)$$

Also Ampère's law [the spatial component of the Yang-

Mills equation (2)] is regained.

$$\partial_t \mathbf{E}_a = i[H, \mathbf{E}_a] = \mathbf{D}_{ab} \times \mathbf{B}_b. \tag{30b}$$

But Gauss' law [the time component of the Yang-Mills equation (2)] is not found.

$$\mathbf{D}_{ab} \cdot \mathbf{E}_b = 0. \tag{30c}$$

(There is no need to seek among the Hamiltonian equations a definition for \mathbf{B} . That quantity is not a fundamental variable, but is given in terms of \mathbf{A} by the conventional formula $\mathbf{B}_a = \nabla \times \mathbf{A}_a - (g/2)\epsilon_{abc}\mathbf{A}_b \times \mathbf{A}_c$.) In a sense, the Hamiltonian theory is a larger theory than the Yang-Mills theory. It gives rise to an entirely consistent quantum mechanics, which, however, does not coincide with the desired gauge theory since Gauss' law is not incorporated.

Let us for the moment ignore this problem and continue the analysis of the Hamiltonian model. It is noted that the theory possesses a symmetry which leaves the equations of motion invariant. In infinitesimal form, the symmetry transformation is

$$\begin{aligned} \delta \mathbf{A}_a &= -g^{-1} \mathbf{D}_{ab} \theta_b, \\ \delta I &= 0, \end{aligned} \tag{31}$$

where θ_b is an arbitrary function of \mathbf{r} , but not of t . Of course we recognize this as just the gauge freedom (7) which respects the $A^0 = 0$ condition, but now we view it as a conventional symmetry which leaves the action invariant. By Noether's theorem we can derive the conserved generator. Since θ_a is an arbitrary local function, one finds local (\mathbf{r} -dependent) conserved quantities.

$$G_a = g^{-1} \mathbf{D}_{ab} \cdot \mathbf{E}_b. \tag{32}$$

(At this stage G_a is nonvanishing since Gauss' law has not as yet been satisfied.) The G_a 's are constants of motion, as is verified by commuting with the Hamiltonian.

$$i[H, G_a] = 0. \tag{33}$$

Gauss' law can now be incorporated into the Hamiltonian quantum theory by demanding that of all the states in the Hilbert space only those that are annihilated by G_a are relevant to the Yang-Mills theory.¹¹

$$G_a |\Psi\rangle = 0. \tag{34}$$

Observe that the G_a 's do not commute.

$$[G_a(\mathbf{r}, t), G_b(\mathbf{r}', t)] = i\epsilon_{abc} G_c(\mathbf{r}, t) \delta(\mathbf{r} - \mathbf{r}'). \tag{35}$$

Hence Eq. (34) is the only possible eigenvalue condition.

An analogy with a simple problem may help in gaining understanding. Consider a particle Hamiltonian for two-dimensional motion, $H = T + V$, $T = \frac{1}{2}p_x^2 + \frac{1}{2}p_y^2$; but with the potential depending only on x , $V = V(x)$. In other words y is an ignorable coordinate in V and the problem possesses a symmetry $\delta y = \alpha$, whose generator is p_y . Suppose further there is a physical requirement on the theory that $p_y = 0$, a condition which

¹¹The earliest use of the present method for quantizing a gauge theory is by Weyl (1950b). Hence the gauge choice $A^0 = 0$ should be called the "Weyl gauge."

can only be imposed on the states $p_y \psi = 0$. This condition is analogous to our Gauss' law Eq. (34); it states that the wave function ψ is independent of the ignorable coordinate y and invariant under the symmetry transformation. Thus similarly Gauss' law (34) requires the Yang-Mills states $\Psi(\mathbf{A})$, viewed as functionals of dynamical variables \mathbf{A}_a , to be independent of those coordinates which are ignorable and which lead to the conservation of G_a . The state must be invariant against transformations generated by G_a .

There is one inconvenient aspect to this development—the states are not normalizable. In the quantum-mechanical example the problem is clear: the integral of $\psi^* \psi$ over y diverges when ψ is y independent. This is a trivial complication which can be removed by simply legislating that the y integration will not be performed, but y will everywhere be set to some preassigned arbitrary value, say 0. Over the remaining variable, x , the normalization integral is still done. (The wave function of course does not depend on y .) The analogous development in the gauge theory will be explained in Sec. V.

Equations (28), (29), and (34) are the basis of a Yang-Mills quantum theory. One may develop from them a perturbative expansion and compute various amplitudes of physical interest. One may also use the above equations to study further properties of the states in theory.

In fact it is convenient, however, to rearrange the perturbation theory, so that it may be represented in terms of conventional diagrams, without constraints like Eq. (34). Indeed in practical computations one uses Faddeev-Popov ghosts, Feynmann Dyson graphs, etc. In Sec. V I, shall show how one passes from the above formulation to the common diagrammatic one. That derivation is useful also for exposing different approaches to the problems of gauge choice in the theory. Before embarking on this route, I wish first to discuss in greater detail the structure of gauge transformations.

IV. GAUGE TRANSFORMATIONS, TOPOLOGY, AND THE VACUUM ANGLE

In this section, I shall remain with the formalism as developed thus far, and study further the action of the gauge symmetry. We have remarked already on the invariance of the quantized theory, when $A^0 = 0$, under transformations which in infinitesimal form are described by Eq. (31), and in finite form by

$$\mathbf{A} \rightarrow U^{-1} \mathbf{A} U - \frac{1}{g} U^{-1} \nabla U. \tag{36a}$$

Here U is a 2×2 unitary, c -number $SU(2)$ matrix, depending on position, but not on time. We shall make a very important hypothesis concerning the physically admissible finite transformations. While some plausible arguments can be given in support of this hypothesis (see below) in the end we must recognize it as an assumption, without which the subsequent development cannot be made. We shall assume that the allowed gauge transformation matrices U tend to a definite limit as r passes to infinity.

$$\lim_{r \rightarrow \infty} U(\mathbf{r}) = U_\infty. \tag{36b}$$

Here U_∞ is a global (position-independent) gauge transformation matrix. With this hypothesis, we are excluding gauge transformations which do not have a well-defined or unique limit at $r \rightarrow \infty$. Equivalently, we are insisting that all physically relevant vector potentials fall faster than $1/r$ at large distances. [By "physically relevant" vector potentials we have in mind those potentials which are arguments of the quantum wave functionals $\Psi(\mathbf{A})$. Alternatively, they may be the potentials over which a functional integration is performed in a function integral formulation.]

Our reasons for adopting this boundary condition are, firstly, the suspicion that potentials not satisfying it are separated by an infinite energy barrier from those that do; hence there will never be any classical transition, and probably no quantal transition either, between them (Jackiw, 1977). Second, the only known physical effects which would require potentials that persist as $1/r$ at large distances involve magnetic monopoles; but these have not yet appeared in Nature, and if they did they would reside in a topologically distinct sector of the Hilbert space, and not in the "vacuum" sector to which the present discussion applies. Third, the total charge

$$Q_a = g \int d\mathbf{r} \varepsilon_{abc} \mathbf{A}_b \cdot \mathbf{E}_c \quad (37a)$$

is well defined only when gauge transformations are restricted. To see this, we use Gauss' law to express Q_a as an integral over a surface at spatial infinity.

$$Q_a = \int d\mathbf{r} \nabla \cdot \mathbf{E}_a = \int dS \cdot \mathbf{E}_a. \quad (37b)$$

Under a local gauge transformation Q_a changes into

$$Q \rightarrow \int dS \cdot U^{-1} \mathbf{E} U \quad (37c)$$

which equals $U_\infty^{-1} Q U_\infty$ when U has a global limit at spatial infinity. Thus for U 's satisfying Eq. (36b), Q transforms by a global gauge transformation. If U has no limit, then the transformed charge has no simple relation to the original charge and we must conclude that the charge is not well defined.

We shall henceforth adopt the hypothesis (36b); the gauge structure of the Yang-Mills theory without this hypothesis has not been thus far determined, though some (isolated) investigations of the monopole sector are available [Witten (1979); Christ and Jackiw (1980); for a review, see Jackiw (1980)].

The insistence on a well-defined limit at spatial infinity is equivalent to compactifying the spatial manifold R^3 to S^3 (the surface of a four-dimensional sphere). The manifold of the $SU(2)$ gauge group is also S^3 , so the matrix functions U provide a mapping $S^3 \rightarrow S^3$ which can be categorized into homotopy classes labeled by an integer, called the "winding number" of the mapping (Jackiw and Rebbi, 1976). In class $n=0$, we place all gauge functions U which are homotopic (continuously deformable) to $U=I$; in class $n=1$, we place those that are not homotopic to I , but rather to some other paradigm, etc. For the $n=1$ paradigm we can choose any gauge function of the form

$$U = e^{i\sigma \cdot \hat{r} f(r)} \quad (38a)$$

with

$$f(\infty) = \pi, \quad f(0) = 0. \quad (38b)$$

This cannot be continuously deformed to I without violating Eq. (36b). The paradigm for class n can be taken to be the above raised to the n th power. Also an analytic expression for the winding number of any gauge transformation is available. It is

$$w = \frac{1}{24\pi^2} \int d\mathbf{r} \varepsilon^{ijk} \text{tr}(U^{-1} \partial_i U)(U^{-1} \partial_j U)(U^{-1} \partial_k U). \quad (39)$$

[Here again we see the need for Eq. (36b); otherwise the integral may diverge.] For gauge functions U satisfying Eq. (36b), w takes on integer values. The existence of gauge transformations which are not homotopic to the identity is what distinguishes the non-Abelian theory from the Abelian one. In the latter, all gauge transformations fall in the homotopically trivial class.

Having classified the gauge functions, we can now inquire how the quantum theory responds to these transformations. It is important to recognize that only gauge transformations which are continuously deformable to the identity, viz. those belonging to the $n=0$ class, can be built up by iterating the infinitesimal transformation generated by G_a ; see also below. Since the infinitesimal generator annihilates the states, the effect of gauge transformations in the trivial homotopy class is to leave the state invariant. Calling the unitary operator which implements the gauge transformation in the n th homotopy class \mathfrak{G}_n , we conclude

$$\mathfrak{G}_0 \Psi(\mathbf{A}) = \Psi(\mathbf{A}). \quad (40)$$

Equation (40) is presented in an explicit Schrödinger picture, where the states are functionals of $\mathbf{A}(\mathbf{r})$. Next we consider the action of \mathfrak{G}_1 . This quantity cannot be obtained by iterating G_a , and we cannot conclude that the state is left invariant. However, since \mathfrak{G}_1 does commute with observables (they are gauge invariant), the only effect it can have on physical states is to leave them phase invariant. We conclude therefore that

$$\mathfrak{G}_1 \Psi(\mathbf{A}) = e^{-i\theta} \Psi(\mathbf{A}) \quad (41)$$

$$\mathfrak{G}_n \Psi(\mathbf{A}) = e^{-in\theta} \Psi(\mathbf{A}).$$

This is the origin of the famous angle in the Yang-Mills theory (Jackiw and Rebbi, 1976; 't Hooft, 1976a,b; Callan, Dashen, and Gross, 1976).

We have no *a priori* computation of θ ; however, one should not entertain the notion that it is an artificial quantity, arising from some peculiarities of our treatment of gauge transformations. Indeed there is another way to see that an arbitrary angle is present in the theory. Let us for a moment suppose that there is no angle in Eq. (41), and that the states are invariant under "large" gauge transformations (those with $n \neq 0$), as well as under "small" gauge transformations (those with $n=0$). We can still find an ambiguity in the choice of the Lagrangian density. Instead of Eq. (4b), one may use

$$\mathcal{L} = -\frac{1}{4} F_a^{\mu\nu} F_{a\mu\nu} + (g^2/32\pi^2) \theta^* F_a^{\mu\nu} F_{a\mu\nu}. \quad (42)$$

The new term is a total divergence; it does not influence the equations of motion.

$$\frac{1}{4} \text{tr} *F^{\mu\nu} F_{\mu\nu} = \partial_\mu X^\mu \tag{43}$$

$$X^\mu = \varepsilon^{\mu\alpha\beta\gamma} \text{tr} \left[\frac{1}{2} A_\alpha \partial_\beta A_\gamma + \frac{1}{3} g A_\alpha A_\beta A_\gamma \right].$$

Consequently, the energy (Hamiltonian) is the same as before, Eq. (28). What is different though, is the relationship between the canonical momentum conjugate to A_a^i and E_a^i . With Eq. (42) we find

$$\Pi_a^i = -E_a^i + (g^2/8\pi^2) \theta B_a^i. \tag{44}$$

Thus the functional Schrödinger equation reads

$$\int dr \left[\frac{1}{2} \left(\frac{\delta}{i \delta A_a^i} - \frac{g^2}{8\pi^2} \theta B_a^i \right)^2 + \frac{1}{2} (B_a^i)^2 \right] \Psi'(A) = E \Psi'(A). \tag{45}$$

The state functionals have been primed to remind that by hypothesis they are invariant under large gauge transformations.

$$\mathfrak{g}_n \Psi'(A) = \Psi'(A). \tag{46}$$

The additional, θ -dependent term in the kinetic energy operator can be eliminated so that the Schrödinger equation acquires the conventional form. To achieve this, we seek a functional of A , $W(A)$, with the property that

$$\frac{\delta W(A)}{\delta A_a^i(r)} = \frac{g^2}{8\pi^2} B_a^i(r). \tag{47}$$

Then the states

$$\Psi(A) = e^{-i\theta W(A)} \Psi'(A) \tag{48}$$

satisfy a Schrödinger equation without θ -dependent terms. Equation (47) may be integrated.

$$W(A) = -\frac{g^2}{4\pi^2} \varepsilon^{ijk} \int dr \text{tr} \left(\frac{1}{2} A^i \partial_j A^k - \frac{1}{3} g A^i A^j A^k \right). \tag{49a}$$

Note that $W(A)$ is essentially the “charge” associated with the “current” X^μ , defined in Eq. (43)

$$W(A) = -\frac{g^2}{4\pi^2} \int dr X^0. \tag{49b}$$

We call $W(A)$ the “winding number” of the gauge potential A .

Next we inquire how $\Psi(A)$ transforms under gauge transformations, given that $\Psi'(A)$ is invariant; clearly that question comes down to determining how $W(A)$ transforms. Upon performing a gauge transformation on $W(A)$, one finds (with the assumed large distance behavior of the relevant quantities)

$$W(U^{-1}AU - g^{-1}U^{-1}\nabla U) = W(A) + w. \tag{50}$$

Here $W(A)$ changes by the winding number of the gauge transformation, and the states $\Psi(A)$, which satisfy a simple Schrödinger equation, are only phase invariant when a large gauge transformation is performed on them.

Thus we see that the angle in the Yang-Mills theory cannot be escaped by simply postulating gauge invariance of the states. It reappears as an ambiguity in the definition of the Lagrangian and of the canonical variables. This alternate point of view has the advantage of demonstrating very clearly that our considerations are gauge invariant— $\theta *F_a^{\mu\nu} F_{a\mu\nu}$ is a gauge-invariant quantity. Also since $*F_a^{\mu\nu} F_{a\mu\nu}$ is CP-odd, the angle is seen to be CP-violating (in theories where

an independent definition of CP can be given).

Finally, let us demonstrate explicitly that a large gauge transformation cannot be obtained by iterating an infinitesimal one. If it were possible, one could write the operator \mathfrak{g} , which implements the finite transformation U , as

$$\mathfrak{g} = \exp -i \int dr g^{-1} (\mathfrak{D}_{ab} \theta_b) \cdot \Pi_a \tag{51a}$$

where

$$U = e^\Omega \tag{51b}$$

$$\Omega = \theta_a \left(\frac{\sigma^a}{2i} \right)$$

and for a large transformation θ_a approaches a nonzero limit at spatial infinity. Let us see whether we can with Eq. (51a) obtain Eq. (50).

$$\mathfrak{g} W(A) = W(A) + \sum_{n=0}^{\infty} \frac{(-1)^{n+1}}{(n+1)!} \left[\int dr g^{-1} (\mathfrak{D}_{ab} \theta_b) \cdot \frac{\delta}{\delta A_a} \right]^n \times \left[\int dr g^{-1} (\mathfrak{D}_{ab} \theta_b) \cdot \frac{\delta}{\delta A_a} \right] W(A). \tag{52a}$$

The last factor above can be evaluated with the help of Eq. (47), and we obtain for it

$$\frac{g}{8\pi^2} \int dr (\mathfrak{D}_{ab} \theta_b) \cdot B_a = \frac{g}{8\pi^2} \int dS \cdot B_a \theta_a - \frac{g}{8\pi^2} \int dr \theta_a \mathfrak{D}_{ab} \cdot B_b. \tag{52b}$$

The last term in Eq. (52b) vanishes by the Bianchi identity [the time component of Eq. (6)]; a result which is just the statement that $G_a W(A) = 0$. In the next-to-last term, the integration is over a sphere at infinity where θ_a attains a nonzero limit. However, in the absence of monopoles, when the potentials fall faster than r^{-1} and the magnetic field falls faster than r^{-2} , the integral over the infinite surface vanishes, leaving

$$\mathfrak{g} W(A) = W(A). \tag{52c}$$

Comparison with Eq. (50) shows that the correct result has been obtained by the above only for gauge transformations with zero winding number. In other words, Eq. (51a) is a correct representation for the operator implementing a gauge transformation belonging to the trivial homotopy class, while gauge transformations in nontrivial homotopy classes cannot be achieved by exponentiating the infinitesimal generator.¹²

The results (41), valid for the Yang-Mills theory, have an analog in a familiar quantum-mechanical problem: particle motion in a periodic potential, e.g. an electron in a crystal. The quantum-mechanical system admits a finite symmetry transformation, translation by the period of the potential. (This is analogous to the large gauge transformation in the Yang-Mills case.) The states are not invariant under this symmetry transformation; rather they are phase invariant, where the

¹²When magnetic monopoles are present and one is working in the monopole sector of the theory where magnetic fields fall as r^{-2} and potentials as r^{-1} at large distances, the situation is obviously different; see Witten (1979), Christ and Jackiw (1980), and Jackiw (1980).

phase involves the Bloch momentum. (The Yang-Mills angle is here the analog.)

Our discussion thus far makes no approximation. But of course, just as in the quantum-mechanical case, a further approximate analysis aids in gaining physical understanding. (In the crystal problem the analogous approach is called the "tight binding" approximation.) One expands the θ -dependent state in a Fourier series.

$$\Psi_\theta = \sum_n e^{in\theta} u_n. \tag{53}$$

The individual u_n 's, which are shifted by a large gauge transformation,

$$g_1 u_n = u_{n+1}, \tag{54}$$

may be considered, in the semiclassical approximation, as describing the system localized in class n of the configuration space. In the crystal these classes are the zones of the periodic potential; in the field theory, they are homotopy classes of the A function space. For sufficiently low energy, there is no classical transition between the individual classes, since an energy barrier separates the classes. However, there are quantum-mechanical, tunneling transitions, and the instantons—imaginary time classical solutions¹³—interpolate in imaginary time between the physically allowed real-time configurations. [It is known that a semiclassical description of tunneling makes use of classical, c -number solutions to dynamical equations, however, not in real time (classically there is no real-time tunneling transition) but in imaginary time.¹⁴ See, for example, Freed (1972) and McLaughlin (1972).] The tunneling picture is certainly useful for a better understanding of the physics of our problem; but it must be emphasized that semiclassical calculations based on instantons and on the instanton gas approximation are reliable only for weak coupling and may not be relevant to the problem at hand when the coupling constant is large. The exact results, based on general considerations sketched above, of course do not rely on any

¹³See Belavin, Polyakov, Schwartz, and Tyupkin (1975), Witten (1977), Chia (1977), 't Hooft (1977), Jackiw, Nohl, and Rebbi (1977), and Atiyah, Drinfeld, Hitchin, and Manin (1978).

¹⁴This is easily seen in the WKB approximation to one-particle quantum mechanics. The real-time classical dynamical equation for a unit-mass particle moving in a potential $V(q)$, $\ddot{q} = -V'(q)$, becomes in imaginary-time [$t \rightarrow -i\tau$] $\ddot{q} = V'(q)$. The first integral gives the imaginary-time energy, which should be zero if we are studying the lowest-lying state. (It is assumed that $V(q)$ has a barrier shape, so that no real-time motion can occur at zero energy.) Hence the equation to solve is $\frac{1}{2}\dot{q}^2 - V(q) = 0$ or $\dot{q} = \pm(2V(q))^{1/2}$ (A). The quantum-mechanical transmission coefficient for penetrating the barrier is determined by e^{-I} , where I is the classical, imaginary-time action evaluated with the solution to (A). $I = \int d\tau (\frac{1}{2}\dot{q}^2 + V(q)) = \int d\tau \dot{q}^2 = \int dq |\dot{q}| = \int dq [2V(q)]^{1/2}$; this is the usual WKB result. Analogously, in the Yang-Mills theory, the imaginary-time energy $\frac{1}{2} \int d\mathbf{r} \{(\mathbf{E}_a)^2 - (\mathbf{B}_a)^2\}$ vanishes for $\mathbf{E}_a = \pm \mathbf{B}_a$. Thus quantum-mechanical tunneling is here described semiclassically by self-dual or anti-self-dual Euclidean [imaginary-time] Yang-Mills fields $F^{\mu\nu} = \pm *F^{\mu\nu}$. By virtue of the Bianchi identity, self-dual and anti-self-dual Yang-Mills fields obviously solve the field equations; they are the celebrated Yang-Mills instantons,¹³ whose further properties are summarized by Jackiw, Nohl, and Rebbi (1978).

semiclassical, weak coupling approximation.

While the quantum-mechanical analogy is a good one, there is an important difference from the Yang-Mills theory. In the crystal, all values of θ are attainable, and θ is a measure of the energy in an energy band. In the Yang-Mills theory, when gauge-invariant quantities are considered, θ cannot change, so that even though one can imagine states with different θ , a complete physical theory is characterized by a single, unique, but as yet undetermined angle. (It is as if in the crystal example, all physical observables were translationally invariant; in that case only one state per band would be physically realizable.)

The coupling of other fields to the Yang-Mills theory does not significantly alter the physical picture presented here, except when the additional fields include massless fermions. In that case, there is a dramatic change. In the presence of massless Fermi fields ψ , there is a new symmetry in the system, chiral invariance, described by the infinitesimal transformation

$$\delta\psi = \gamma_5 \psi. \tag{55a}$$

Noether's theorem gives rise to the axial-vector current.

$$J_5^\mu = i\bar{\psi} \gamma^\mu \gamma^5 \psi. \tag{55b}$$

However, as is well known, quantal renormalization effects interfere with the conservation of this current. Rather than being conserved, J_5^μ is afflicted by the axial-vector current anomaly.⁸

$$\partial_\mu J_5^\mu = g^2 \frac{c}{8\pi^2} \text{tr} *F^{\mu\nu} F_{\mu\nu}. \tag{56}$$

The constant c depends on the number of fermion species, while the matrices $F^{\mu\nu}$ over which the trace is taken refer to the representation matrices of the fermions. For one species of fermions in the fundamental (isospin 1/2) representation, $c=1$ and $F^{\mu\nu}$ is constructed from the Pauli matrices as above. From Eq. (43) we see that Eq. (56) implies that a conserved current is

$$J_5^\mu = j_5^\mu - \frac{g^2}{2\pi^2} X^\mu \tag{57a}$$

and the conserved charge involves the quantity already encountered in Eq. (49).

$$Q_5 = q_5 + 2W(\mathbf{A}). \tag{57b}$$

The total, time-independent chiral charge Q_5 has two pieces: q_5 , which arises from the fermions and is gauge invariant; $2W(\mathbf{A})$ the gauge field contribution, which according to Eq. (50) is not gauge invariant. This means that also Q_5 is not gauge invariant against large gauge transformations; it changes by twice the winding number.

There are now three operators to consider: H , g_n , and Q_5 . The Hamiltonian commutes with the other two, but they do not commute with each other.

$$[g_n, Q_5] = 2n g_n. \tag{58}$$

The three cannot be simultaneously diagonalized. Gauge invariance requires that g_n be diagonalized, as in Eq. (41), but that means that Q_5 acts as a raising operator and changes θ

$$e^{-i\theta' Q_5} \Psi_\theta = \Psi_{\theta+2\theta'} \quad (59)$$

The energy eigenvalues of H , which commutes with Q_5 , can no longer depend on θ ; tunneling is suppressed and the entire energy band collapses to one level. Moreover, the chiral symmetry is spontaneously broken since the states are not chirally invariant. The calculational, practical consequences of all this are limited by the approximations that are used, but again we see that certain exact results can be established.

While the main thrust of this review is to present a development of the Yang-Mills quantum theory, and not to discuss the role that it plays in phenomenological applications to model building for strong chromodynamics or electroweak flavor dynamics, the topological results presented in this section are sufficiently novel and subtle that it is quite appropriate here to remark on their physical significance. In fact, the unexpected intrusion of topology into quantum physics served to resolve a long-standing conflict between the (apparent) predictions of the Yang-Mills theory as applied to particle phenomenology and physical reality. The puzzle, called the $U(1)$ problem, was the following. [For a review, see Weinberg (1975).] The quantum chromodynamical Lagrangian, posited as governing the dynamics, is constructed from Yang-Mills fields and massless Fermi fields. Consequently it possesses a high degree of chiral symmetry which is not observed in Nature, and theorists have supposed that the observed particle spectrum reflects the fact that this large symmetry is not realized in Wigner–Weyl multiplets, but in the Nambu–Goldstone fashion based on spontaneous symmetry breakdown. In this way one avoided most of the undesirable consequences of the large symmetry, but one particular symmetry transformation remained unaffected, giving unwanted and unobserved predictions. The symmetry, a universal chiral rotation of all the massless Fermi fields in the theory, described by Eq. (55a), appears to give rise to a conserved, singlet axial vector current, Eq. (55b), which in a Wigner–Weyl realization predicts the conservation of the number of right-handed fermions minus that of the left-handed ones. This regularity is not observed, but the Nambu–Goldstone realization also leads to an unacceptable result: there should be a massless Goldstone boson degenerate with the pions. (In the approximation that we are using the pions are massless.) Unfortunately, no such particle is observed.

The first indication that something unexpected may be transpiring came with the discovery of the anomaly (56) in the axial-vector current.⁸ Yet this did not immediately resolve the difficulty since the anomaly appeared too “soft” to remove the unwanted predictions. Specifically, if one computed the total change in fermion chirality $q_5(t) = \int d\mathbf{r} j_5^0(x)$

$$\begin{aligned} \Delta q_5 &= q_5(\infty) - q_5(-\infty) = \int_{-\infty}^{\infty} dt \frac{d}{dt} q_5(t) = \int d^4x \partial_\mu j_5^\mu \\ &\propto \int d^4x \operatorname{tr} {}^*F^{\mu\nu} F_{\mu\nu} \end{aligned} \quad (60)$$

it did not seem that the anomaly gave rise to chiral

nonconservation since ${}^*F^{\mu\nu} F_{\mu\nu}$ is a total divergence and one believed that surface terms can be ignored.

The next step was the discovery by Belavin, Polyakov, Schwartz, and Tyupkin (1975) of self-dual Yang-Mills configurations in Euclidean space, called “instantons” or “pseudoparticles,” for which the integral of $\operatorname{tr} {}^*F^{\mu\nu} F_{\mu\nu}$ is nonvanishing. Indeed, they drew the attention of the physics community to the fact that the integral of ${}^*F^{\mu\nu} F_{\mu\nu}$, known as the “Pontryagin Index” to mathematicians, is a measure of the topological properties of gauge fields; its nonzero value is quantized. Then ‘t Hooft (1976a,b) suggested that these configurations be used in an approximate, semiclassical evaluation of the functional integral, continued to Euclidean space. He calculated various chirality-changing amplitudes, and showed that in this approximation they were nonvanishing, indicating an absence of the unwanted chiral symmetry. In the notation of Eq. (60), Δq_5 does not vanish, since the integral is nonzero for instantons. This provided a resolution of the $U(1)$ problem, a resolution which may now be offered independently of the semiclassical approximation, when it is realized that the vacuum structure of the theory prevents states from possessing the chiral symmetry; see Eq. (59).

While the general considerations remove the $U(1)$ problem, the calculational, practical consequences of the nontrivial topology are limited by the approximations that are used. Specifically, the semiclassical results are numerically reliable only for weak coupling, and as yet it is not clear whether the physical theory can in fact be described by such a small coupling constant.

Finally, mention need be made of the angle θ . Since it is a CP -violating parameter, the experimental limits on the value of the neutron electric dipole moment put strong constraints on its value: it must be vanishingly small.¹⁵ This presents something of a conceptual problem since at present we do not know what determines the value of θ , and what mechanism assures its negligible magnitude. It should be realized that one cannot set it zero *ab initio*. The reason is the following. When we consider massive quarks, as we must, the mass arising from spontaneous symmetry breaking and from radiative corrections appears in the Lagrangian in the form $C_1 \bar{\psi}\psi + C_2 \bar{\psi}\gamma^5\psi$. In order to isolate the CP -violating portions of the theory, one needs to remove the apparent P violating mass term involving C_2 . This can be done by a chiral rotation, which formally leaves invariant the remainder of the Lagrangian (kinetic terms for the quarks, Yang-Mills terms, quark–gauge field interaction terms). However, because of the axial-vector anomaly, the chiral rotation produces a contribution in the Lagrangian proportional to $\operatorname{tr} {}^*F^{\mu\nu} F_{\mu\nu}$ (Gross and Jackiw, 1972), and comparison with Eq. (42) shows that even if θ is initially set to zero, a new angle emerges when the mass is rediagonalized. In other words, in order to have physically vanishing θ , one must fine-tune a “bare” angle so it precisely cancels the effects

¹⁵Crewther, DiVecchia, Veneziano, and Witten (1979) find $\theta \lesssim 10^{-9}$.

of a mass rediagonalization. The mechanism for this fine tuning is at present unknown.¹⁶

V. SOLVING GAUSS' LAW CONSTRAINTS

We return now to the main line of development of the canonical theory. Our purpose here is to show how one can eliminate the ignorable variables in the problem, viz. we set to zero the momentum conjugate to the ignorable coordinate, thus explicitly satisfying Gauss' law, and we evaluate the ignorable coordinate at some definite value, thus removing the divergence in the normalization integral. These calculations are most easily done with the help of the functional integral, and we begin by writing that integral in Hamiltonian form.¹⁷

$$Z = \int \{d\mathbf{E}_a\} \{d\mathbf{A}_a\} \delta(G_a) \delta(\chi_b) \\ \times \exp - i \int d^4x [\mathbf{E}_a \cdot \partial_t \mathbf{A}_a + \frac{1}{2} (\mathbf{E}_a)^2 + \frac{1}{2} (\mathbf{B}_a)^2].$$

Here the exponent involves the field-theoretical generalization of $\int dt [p\dot{q} - H]$; we take the canonical momentum to be the negative of the electric field. The first delta function enforces Gauss' law, viz. the vanishing of the components of G ; the second evaluates the ignorable coordinates conjugate to G at some pre-assigned value, in order to remove the non-normalizability of states. The ignorable coordinates have here been designated by χ , quantities which are to a large extent arbitrary. Indeed, they need not even be conjugate to G ; provided the Poisson bracket $\{\chi_a, \chi_b\}$ vanishes, we may use χ 's for which the Poisson bracket $\{G_a, \chi_b\}$ is not the identity, but then an additional factor appears in the functional integral. In this more general case we use (Faddeev, 1970)

$$Z = \int \{d\mathbf{E}_a\} \{d\mathbf{A}_a\} \delta(G_a) \delta(\chi_b) \det\{G_a, \chi_b\} \\ \times \exp - i \int d^4x [\mathbf{E}_a \cdot \partial_t \mathbf{A}_a + \frac{1}{2} (\mathbf{E}_a)^2 + \frac{1}{2} (\mathbf{B}_a)^2], \quad (61)$$

where the functional determinant compensates for any noncanonical structure in χ . The value of the functional integral can be shown to be independent of the choice of χ (Faddeev, 1970).

Let us now convert the above Hamiltonian formulation into the more familiar Faddeev-Popov Lagrangian theory (Faddeev and Popov, 1967). To do this, we assume that χ depends only on the coordinates \mathbf{A} , but not on the momenta \mathbf{E} . Since G is the generator of infinitesimal gauge transformations, $\{G_a, \chi_b\}$ is nothing but the infinitesimal gauge transform of χ , hence

¹⁶The above remarks discuss the physical role of topology in quantum chromodynamics. 't Hooft (1976a, b) pointed out that also quantum flavor dynamics is affected, by topological non-conservation of conventional quantum numbers. This may induce for example proton decay, but with exceedingly small probability, hence it is practically irrelevant.

¹⁷We shall use the functional integral heuristically, not paying attention to questions of well definition. Consequently the corresponding questions of operator ordering will be here ignored. By using operator methods or by careful analysis of the functional integral one may of course regain these terms. A recent study of this topic for Yang-Mills theory is by Christ and Lee (1980).

also independent of \mathbf{E} . In an obvious notation, we may designate that Poisson bracket by $\delta_a \chi_b$. Thus the \mathbf{E} dependence resides quadratically in the exponential and also in $\delta(G_a)$, allowing for an evaluation of the functional \mathbf{E} integral. This is achieved by first writing

$$\delta(G_a) = \int \{dA_a^0\} \exp i \int d^4x A_a^0 G_a. \quad (62)$$

Then the \mathbf{E} integral in Eq. (61) is Gaussian; after it is performed one is left with

$$Z = \int \{dA_a^\mu\} \delta(\chi_a) \det(\delta_a \chi_b) \exp i \int d^4x \mathcal{L} \quad (63)$$

with $\mathcal{L} = \frac{1}{2} (\mathbf{E}_a^2 - \mathbf{B}_a^2)$ and \mathbf{E} no longer an independent variable, but given in terms of the potential, $\mathbf{E} = -\partial^0 \mathbf{A} - \nabla \mathbf{A}^0 + g[\mathbf{A}, \mathbf{A}^0]$. Equation (63) is recognized on the familiar Faddeev-Popov expression: the delta function is the "gauge choice," the determinant is the "gauge compensating" factor which may also be represented by a Faddeev-Popov ghost integral (Faddeev and Popov, 1967).

There are two reasons why I took so many steps to arrive at the Faddeev-Popov formula, which usually is "derived" by starting from the Lagrangian functional integral $\int \{dA_a^\mu\} \exp i \int d^4x \mathcal{L}$ and "canceling" an infinite factor coming from the integration over the group. Firstly, I wanted to show how Eq. (63) follows from the Hamiltonian formalism, which is the only one that is entirely reliable. When one begins with the Lagrangian formulation, which is obviously more elegant and compact, mistakes can occasionally be made, which are rectified only when it is realized that the (correct) Hamiltonian formulation does not imply the naive (incorrect) Lagrangian formulation in the general case. Rather a modification must be made. [The most recent example of this recurring phenomenon had to do with the 4-ghost interaction in supergravity, which was initially missed in the Lagrangian formulation.] A second reason for beginning with the canonical Hamiltonian is that it sometimes happens that the integration over \mathbf{E} is not as straightforwardly carried out as in the above example and a Lagrangian formulation cannot even be attained; see below.

The next subject to discuss is the choice of χ . In the Abelian, Maxwell case the most natural choice, and the one most frequently made, is $\chi = \nabla \cdot \mathbf{A}$, i.e., one sets to zero the longitudinal vector potential. This "Coulomb-gauge" choice is indeed appropriate since the longitudinal component of the vector potential is the ignorable coordinate: \mathbf{B}^2 does not depend on it. Also $\nabla \cdot \mathbf{A}$ is virtually canonically conjugate to the Abelian Gauss' law generator $\nabla \cdot \mathbf{E}$: their Poisson bracket is just the Laplacian, a constant quantity which may be ignored in the functional integral. In this way one arrives from Eqs. (61) and (63) at the usual, Coulomb-gauge, quantization of the Maxwell theory.

When Yang-Mills theory was first considered, the Coulomb gauge was again used for quantization (Schwinger, 1962); the analogy was drawn with electromagnetism, but the question of whether this was an appropriate and natural choice was not addressed. It is, however, clear, that unlike in the Abelian model, the longitudinal vector potentials are not ignorable co-

ordinates; $\nabla \cdot \mathbf{A}$ does occur in \mathbf{B} through the terms quadratic in the potential. Also $\nabla \cdot \mathbf{A}$ is far from being conjugate to G ; $\delta_a \chi_b$ is $\mathfrak{D}_{ab} \cdot \nabla$ which depends on \mathbf{A} . There are additional difficulties with this choice of χ since $\det \mathfrak{D}_{ab} \cdot \nabla$ can vanish for some (large) values of \mathbf{A} and the whole procedure becomes ill-defined.¹⁸ In spite of these shortcomings, the Coulomb gauge (or its various generalizations and modifications) is widely used in perturbative calculations for the Yang-Mills theory. Of course perturbation theory is an expansion around an Abelian, Maxwell-like limit and for this reason one can make a case for this gauge choice. Also in perturbation theory one never sees the zeroes of $\mathfrak{D}_{ab} \cdot \nabla$. With the Coulomb gauge choice for χ , Eq. (63) provides the conventional Faddeev-Popov quantization of the Yang-Mills theory. It is these perturbative calculations that have given us all the evidence for the relevance of Yang-Mills theory to a description of physical processes.

However, one wants to look beyond perturbation theory. An interesting problem is to find the proper ignorable coordinates and to construct an effective Hamiltonian in terms of unconstrained variables. This problem has been solved both by direct quantum-mechanical methods (Goldstone and Jackiw, 1978; and Baluni and Grossman, 1978) and by functional integral techniques (Izergin, Korepin, Semenov-Tian-Shansky, and Faddeev, 1979). A convenient choice for the χ 's expresses them in terms of the E 's and not in terms of the A 's. Although the effective Hamiltonian can then be given explicitly, precisely because the approach is nonperturbative, no practical applications of the formalism have thus far been made. The Hamiltonian contains inverse powers of the coupling constant which frustrate naive approximation methods. We present here only a brief outline of these ideas, as realized by functional methods. Those interested in more details are referred to the cited literature [see Goldstone and Jackiw (1978), Baluni and Grossman (1978), and Izergin *et al.* (1979)].

In Eq. (61) χ is chosen as follows

$$\chi_a = \varepsilon_{aib} E_b^i \tag{64a}$$

i.e., the three ignorable coordinates are taken to be the antisymmetric portion of E_a^i , viewed as a 3×3 matrix in the combined space, isospace components. It follows that the gauge compensating term is $\det \mathfrak{M}_{ab}$

$$\mathfrak{M}_{ab}(\mathbf{x}, \mathbf{y}) = (E_b^a - \delta_{ab} E_c^c) \delta(\mathbf{x} - \mathbf{y}). \tag{64b}$$

In order to find the Hamiltonian for the unconstrained variables, which are $E_{(ia)}$, the symmetric portion of the 3×3 matrix E_a^i , and $A_{(ia)}$, the symmetric portion of

the 3×3 matrix A_a^i , we perform the integration in Eq. (61) over the corresponding antisymmetric parts. The electric field integral is trivial, since the delta function simply sets the antisymmetric part to zero, leaving

$$Z = \int \{dE_{(ia)}\} \{dA_a^i\} \delta(G_a) \det(\mathfrak{M}_{ab}) \times \exp - i \int d^4x [E_{(ia)} \partial_t A_{(ia)} + \frac{1}{2} (E_{(ia)})^2 + \frac{1}{2} (B_a^i)^2]. \tag{65a}$$

Next we decompose A_a^i into its symmetric and antisymmetric parts and use the remaining delta function to integrate over the latter.

$$A_a^i = A_{(ia)} + \varepsilon_{iab} A^b \tag{65b}$$

$$\begin{aligned} \delta(G_a) \det(\mathfrak{M}_{ab}) &= \delta(\partial_i E_{(ia)} - g \varepsilon_{abc} A_{(ib)} E_{(ic)} - g \mathfrak{M}_{ab} A^b) \det(\mathfrak{M}_{ab}) \\ &= \delta(g^{-1} \mathfrak{M}_{ab}^{-1} (\partial_i E_{(ia)} - g \varepsilon_{abc} A_{(ib)} E_{(ic)}) - A^a). \end{aligned} \tag{65c}$$

We see from Eq. (65c) that the product of the Gauss' law delta function with the gauge compensating determinant yields a delta function which evaluates A^a , the antisymmetric part of A_a^i , without any further determinant. Thus we find from Eq. (65a)

$$Z = \int \{dE_{(ia)}\} \{dA_{(ia)}\} \times \exp - i \int d^4x [E_{(ia)} \partial_t A_{(ia)} + \frac{1}{2} (E_{(ia)})^2 + \frac{1}{2} (B_a^i)^2]. \tag{65d}$$

This shows that the unconstrained canonically conjugate variables are $A_{(ia)}$ and $-E_{(ia)}$, while the Hamiltonian governing dynamics is

$$H = \frac{1}{2} \int d^3r \{ (E_{(ia)})^2 + (B_a^i)^2 \}, \tag{66}$$

where it is understood that B_a^i is constructed in the conventional way from the full A_a^i , whose antisymmetric part $\varepsilon_{iab} A^b$ is a dependent quantity, given in terms of the canonical variable by

$$A^b = g^{-1} \mathfrak{M}_{bc}^{-1} (\partial_i E_{(ic)} - g \varepsilon_{cde} A_{(id)} E_{(ie)}). \tag{67}$$

The Hamiltonian is singular at those points in $E_{(ic)}$ function space where \mathfrak{M}_{bc} has no inverse. These are of course the same points where $\det \mathfrak{M}_{ab}$ vanishes. It does not seem that this singularity has any direct dynamical significance. Rather it is of kinematical origin, arising from the fact that the coordinates in function space that we are using cannot be defined on the entire space without singularities. Analogy with radial coordinates in a one-particle quantum-mechanical problem may be drawn. There too the origin provides a kinematical singularity where the angles cannot be defined. Correspondingly, the Hamiltonian is singular at that point (centrifugal barrier), but these effects do not have a dynamical origin.

Note the occurrence of the coupling constant in the

¹⁸See Mandelstam (1977) and Gribov (1977, 1978). These difficulties with the Coulomb condition are a consequence of the nontrivial topology carried by the Yang-Mills potentials; see Jackiw, Muzinich, and Rebbi (1978) and Ademollo, Napolitano, and Sciuto (1978). Indeed, even in the Maxwell theory, the Coulomb gauge condition becomes beset by intricacies in the presence of topologically nontrivial structures, viz. magnetic monopoles. That these problems arise also when any other gauge condition is imposed on \mathbf{A} has been shown, under additional hypothesis, by Singer (1978). For a summary, see Jackiw (1978).

demoninator. The limit $g=0$ is no longer attainable; this is a reflection of the fact that we have integrated out completely and exactly the non-Abelian gauge sector of the theory, which necessarily involves non-vanishing g . (When $g=0$, the theory does not possess a non-Abelian invariance; rather it supports a direct product of Abelian symmetry groups.) While it is suggestive to speculate that this formalism is relevant to strong coupling investigations, it has not been possible, as yet, to effect any useful calculations. Nevertheless it appears that our gauge choice will be employed in approximation schemes which do not rely on expansions in the coupling constant. For example, it has been used in large N analyses of $SU(N)$ gauge theories by Baluni (1980).

VI. CONCLUSION

The development here has been formal, but based on elementary, familiar principles of quantum mechanics. In this way we have reached by canonical quantization methods the usual Faddeev-Popov formulation, Eq. (63), which is then the starting point of perturbative investigations. Because of asymptotic freedom these have had successful application to high-energy phenomenology in quantum chromodynamics. Also because of the smallness of electroweak couplings, perturbative calculations have correctly described the low-energy regime of quantum flavor dynamics. Within the canonical framework, we exposed in Sec. IV some of the unexpected richness in the theory which is a consequence of nontrivial topology. Here numerical calculations are based on the semiclassical technique whose reliability is uncertain, hence it is good to have some exact results which emerge without approximation. Finally, in the second half of Sec. V, unconventional approaches to the quantum theory were suggested. These have not as yet produced further illumination of the model, though Baluni's (1980) work is an interesting attempt in that direction.

ACKNOWLEDGMENTS

I wish to thank E. Calva-Tellez, D. V. and G. V. Chudnovsky, as well as S. Hawking for giving me the opportunity to present this material to an interested and stimulating audience. The research was supported in part through funds provided by the U. S. Department of Energy (DOE) under contract EY-76-C-02-3069.

REFERENCES

- Ademollo, M., E. Napolitano, and S. Sciuto, 1978, Nucl. Phys. B **134**, 477.
- Atiyah, M., V. Drinfeld, N. Hitchin, and Y. Manin, 1978, Phys. Lett A **65**, 185.
- Baluni, V., 1980, Phys. Lett. B **90**, 407.
- Baluni, V. and B. Grossman, 1978, Phys. Lett. B **78**, 226.
- Belavin, A., A. Polyakov, A. Schwartz, and Y. Tyupkin, 1975, Phys. Lett. B **59**, 85.
- Bergmann, P., and E. Flaherty, 1978, J. Math. Phys. **19**, 212.
- Callan, C., R. Dashen, and D. Gross, 1976, Phys. Lett. B **63**, 334.
- Chia Kwei Peng, 1977, Sci. Sin. **20**, 345.
- Christ, N., and R. Jackiw, 1980, Phys. Lett. B **91**, 228.
- Christ, N., and T. D. Lee, 1980, Phys. Rev. D **22**, 939.
- Coleman, S., 1977, in *New Phenomena in Subnuclear Physics*, edited by A. Zichichi (Plenum, New York).
- Coleman, S., 1979, in *The Phys of Subnuclear Physics*, edited by A. Zichichi (Plenum, New York).
- Crewther, R., P. DiVecchia, G. Veneziano, and E. Witten, 1979, Phys. Lett. B **88**, 123.
- Deser, S., M. Duff, and C. Isham, 1980, Phys. Lett. B **93**, 419.
- Faddeev, L., 1970, Teor. Mat. Fiz. **1**, 1 [Theor. Math. Phys. (USSR) **1**, 1 (1970)].
- Faddeev, L., and V. Popov, 1967, Phys. Lett. B **25**, 29.
- Faddeev, L. and A. Slavnov, 1980, *Gauge Fields* (Benjamin, Reading MA).
- Fairlie, D., 1979, Phys. Lett. B **82**, 97.
- Forgács, P., and N. Manton, 1980, Commun. Math. Phys. **72**, 15.
- Freed, K., 1972, J. Chem. Phys. **56**, 692.
- Friedman, J. and R. Sorkin, 1980, Phys. Rev. Lett. **44**, 1100.
- Goldstone, J., and R. Jackiw, 1978, Phys. Lett. B **74**, 81.
- Gribov, V., 1977, Lecture at 12th Winter School, Leningrad.
- Gribov, V., 1978, Nucl. Phys. B **139**, 1.
- Gross, D., and R. Jackiw, 1972, Phys. Rev. D **6**, 477.
- Gross, D., and F. Wilczek, 1973, Phys. Rev. Lett. **30**, 1343.
- Harnad, J., S. Shnider, and L. Vinet, 1980, in *Mathematical Problems in Theoretical Physics*, edited by K. Osterwalder (Springer, Berlin).
- Isham, C., 1980, to appear in the Wolfgang Youngrau memorial volume [Imperial College preprint ICTP/79-80/45].
- Izergin, A., V. Korepin, M. Semenov-Tian-Shansky, and L. Faddeev, 1979, Teor. Mat. Fiz. **38**, 3 [Theor. Math. Phys. (USSR) **38**, 1 (1979)].
- Jackiw, R., 1972, in *Lectures on Current Algebra and Its Applications*, by S. Treiman, R. Jackiw, and D. Gross (Princeton University, Princeton, NJ).
- Jackiw, R., 1977, Rev. Mod. Phys. **49**, 681.
- Jackiw, R., 1978, in *New Frontiers in High-Energy Physics*, edited by B. Kursunoglu, A. Perlmutter, and L. Scott (Plenum, New York).
- Jackiw, R., 1979, Phys. Rev. Lett. **41**, 1635.
- Jackiw, R., 1980, Acta Phys. Austriaca (in press) [Institute for Theoretical Physics preprint, NSF ITP 80-15].
- Jackiw, R., and N. Manton, 1980, Ann. Phys. (NY) **127**, 257.
- Jackiw, R., I. Muzinich, and C. Rebbi, 1978, Phys. Rev. D **17**, 1576.
- Jackiw, R., C. Nohl, and C. Rebbi, 1977, Phys. Rev. D **15**, 1642.
- Jackiw, R., C. Nohl, and C. Rebbi, 1978, in *Particles and Fields*, edited by D. Boal and A. Kamal (Plenum, New York).
- Jackiw, R., and C. Rebbi, 1976, Phys. Rev. Lett. **37**, 172.
- Kibble, T., 1961, J. Math. Phys. **2**, 212.
- Klein, O., 1939, in *New Theories in Physics* (International Institute of Intellectual Cooperation).
- Kobayashi, S., and K. Nomizu, 1963, *Foundations of Differential Geometry* (Interscience, New York).
- Mandelstam, S., 1977, Lecture at American Physical Society Meeting, Washington, D. C.
- Manton, N., 1979, Nucl. Phys. B **158**, 141.
- Marciano, W., and H. Pagels, 1978, Phys. Rep. C **36**, 137.
- Mayer, M., 1980, in *Mathematical Problems in Theoretical Physics*, edited by K. Osterwalder (Springer, Berlin).
- McLaughlin, D., 1972, J. Math. Phys. **13**, 1099.
- Politzer, H., 1973, Phys. Rev. Lett. **30**, 1346.
- Sakurai, J. J., 1960, Ann. Phys. (NY) **11**, 1.
- Salam, A., 1968, in *Relativistic Groups and Analyticity*, edited by N. Svartholm (Interscience, New York).
- Schwarz, A., 1977, Commun. Math. Phys. **56**, 79.
- Shaw, R., 1955, Ph.D. Thesis, Cambridge University.
- Schwinger, J., 1962, Phys. Rev. **125**, 1043.
- Sciama, D., 1962, *Recent Developments in General Relativity* (Pergamon, Oxford), p. 415.
- Singer, I., 1978, Commun. Math. Phys. **60**, 7.
- Taylor, J., 1976, *Gauge Theories of Weak Interactions*, (Cambridge University, Cambridge, England).

- 't Hooft, G., 1972, unpublished.
- 't Hooft, G., 1976a, Phys. Rev. Lett. **37**, 8.
- 't Hooft, G., 1976b, Phys. Rev. D **14**, 3432.
- 't Hooft, G., 1977, unpublished.
- Trautman, A., 1979, Bull. Acad. Pol., Ser. Sci. Phys. Astron. Phys. **27**, 7.
- Weinberg, S., 1967, Phys. Rev. Lett. **19**, 1264.
- Weinberg, S., 1975, Phys. Rev. D **11**, 3583.
- Weyl, H., 1950a, *Space, Time, Matter* (Dover, New York).
- Weyl, H., 1950b, *The Theory of Groups and Quantum Mechanics* (Dover, New York).
- Witten, E., 1977, Phys. Rev. Lett. **38**, 121.
- Witten, E., 1979, Phys. Lett. B **86**, 283.
- Wu, T. T., and C. N. Yang, 1975, Phys. Rev. D **12**, 3845.
- Yang, C. N., and R. Mills, 1954, Phys. Rev. **96**, 191.