

Large gauge transformations and the strong CP problem*

John Dougherty

October 2, 2019

Abstract

According to the Standard Model of particle physics, some gauge transformations are physical symmetries. That is, they are mathematical transformations that relate representatives of distinct physical states of affairs. This is at odds with the standard philosophical position according to which gauge transformations are an eliminable redundancy in a gauge theory’s representational framework. In this paper I defend the Standard Model’s treatment of gauge from an objection due to Richard Healey. If we follow the Standard Model in taking some gauge transformations to be physical symmetries then we face the “strong CP problem”, but if we adopt the standard philosophical position on gauge then the strong CP problem dissolves. Healey offers this as a reason in favor of the standard philosophical view. However, as I argue here, following Healey’s recommendation gives a theory that makes bad empirical predictions.

1 Introduction

According to philosophical orthodoxy, gauge transformations signal a “redundant descriptive apparatus. . . a veil through which the intrinsic gauge-independent content of the model can be dimly glimpsed” (Earman, 2004, 189). They represent “an unphysical symmetry which merely relates different representations of the same physical state”, and they “could in principle be eliminated by passing to a reformulation in terms of gauge-invariant degrees of freedom” (Struyve, 2011, 226). You might say they’re not symmetries at all: they are redundancies, and a “more intrinsic” formulation of gauge theories would feature one representative for every gauge equivalence class of gauge configurations, hence “would neither be, nor fail to be gauge symmetric” (Healey, 2007, 185). In short, gauge transformations are an artefact of the theory’s formalism. If we want to understand the physical content of the theory, it would be better to do away with them.

This eliminative interpretation of gauge is at odds with the our current best theory of high-energy physics. In the Standard Model of particle physics—and particularly in quantum chromodynamics (QCD), the model of the strong force—gauge transformations are only said to be redundancies if they act as

*This is a preprint of an article appearing in *Studies in the History and Philosophy of Modern Physics* (doi:10.1016/j.shpsb.2019.09.001). Please only cite the final version.

the identity at infinity. A “small” gauge transformation of this kind relates different representations of the same state of affairs, in line with the sentiments of the opening paragraph. But a “large” gauge transformation—one which does not act as the identity at infinity—is a real physical symmetry. That is, a large gauge transformation relates representatives of different physical states. Mathematical differences between these representatives can reflect a physical difference, signaling the existence of some quantities and possibilities that cannot exist according to the received philosophical position.

In this paper I defend the physical significance of the distinction between large and small gauge transformations against the eliminative interpretation of gauge. In particular, I respond to an argument from Richard Healey’s (2007, 2010) sustained defense of the triviality of gauge structure. Healey argues that the Standard Model is afflicted with an avoidable explanatory problem. If configurations related by a large gauge transformation are distinct then mathematical differences between them can represent physical differences. These potential physical differences pick out a preferred direction in time. But the observed physics of the strong force is invariant under time reversal—equivalently, it is invariant under the CP operation, the combination of charge conjugation and parity inversion. So the Standard Model faces a “strong CP” problem: why does the physics of the strong force appear to be CP-invariant, given that CP-violating processes are possible?

Healey argues that this strong CP problem is a symptom of the Standard Model’s confusion about gauge. On the eliminative view, “all states and variables are automatically invariant under both ‘small’ and ‘large’ gauge transformations” (Healey, 2007, 198). They must be, since these transformations are a feature of the mathematics of the theory and not a feature of the world. As Healey puts it, a “more intrinsic formulation of a classical Yang–Mills theory would not even mention gauge, and so the issue of its gauge symmetry would not arise” (2007, 185). If this were true, then the mathematical differences that lead to CP violation in the Standard Model should be interpreted as merely mathematical differences that reflect no physical features. CP violation in strong force physics is thus shown to be impossible. And not a moment too soon: expenditure of time and money on the experimental search for this violation has been steadily increasing (Irastorza and Redondo, 2018)!

In the following I lay out the empirical case for the size distinction in response to Healey’s argument from the strong CP problem. A theory that treats large gauge transformations as redundancies makes wrong predictions. There are two ways to understand Healey’s proposed solution to the strong CP problem. The first reading gives a theory in which particles live one hundred times longer than they do in reality. The second predicts an undiscovered particle that we would have seen in synchrotrons and cosmic rays in the early 1950s, had it existed. Neither of these predictions is acceptable. Indeed, these problems were once thought to be fatal threats to QCD and its precursors. The distinction between large and small gauge transformations appears in the Standard Model today because ’t Hooft (1976a,b, 1986) used it to show that QCD doesn’t predict this nonexistent particle, building on Adler (1969) and Bell and Jackiw’s (1969) solution to the problem of particle decay widths.

I will conclude that an eliminative view of gauge is empirically inadequate and I will suppose that Healey’s response to the strong CP problem consequently fails. I will not attempt to weigh the costs of the strong CP problem against

the costs of empirical inadequacy, nor will I try to pin down what exactly is problematic about it. Nor do the derivations of empirical predictions that I will discuss depend on the strong CP problem: live solutions to the strong CP problem are also compatible with these predictions. So the strong CP problem is ancillary to the main thrust of my argument. It plays two supporting roles. First, I will use the theoretical context of the strong CP problem to fill in the details of Healey’s objection. Healey dissolves the strong CP problem by rejecting one feature of the Standard Model but does not describe a positive alternative. To show that the eliminative view falls afoul of observation, I will argue that modifying the Standard Model as Healey suggests can only avoid the strong CP problem at the cost of empirical adequacy. The presentation of the strong CP problem in Section 2 therefore serves as a fiducial theory from which Healey’s theory is a departure. Second, the strong CP problem serves to connect the results I will discuss with the eliminative view—while proponents of the eliminative view have not dealt directly with these empirical predictions, I take Healey’s arguments to show that the view must reject the presuppositions of the strong CP problem. My argument in this paper aims to show that the eliminative view is thereby committed to falsified predictions.

I begin with a relatively self-contained statement of the strong CP problem and its context. This is a well-known problem in QCD that appears in most textbooks on quantum field theory, but it has a fair few ingredients and tends to be mischaracterized in the philosophical literature. In particular, most philosophical discussions of the strong CP problem focus exclusively on the strong force itself, ignoring the matter fields. But as I explain in Section 2, CP violation in the matter sector can cancel CP violation by the strong force, and the strong CP problem is concerned with whatever CP violation remains after this cancellation.¹ In Section 3 I argue that the disagreement over large gauge transformations amounts to a disagreement over the value of a particular term in the theory’s action. This term varies under large gauge transformations, so if large gauge transformations are redundancies then the term is ill-defined. If large gauge transformations are merely symmetries with respect to some properties then this term can be nonzero and can lead to CP violation in the strong sector. Framing the disagreement in this way is consistent with Healey’s remarks and gives two ways of spelling out his objection in concrete terms. In Section 4 I show that both of these options lead to bad predictions.

2 The strong CP problem

The strong CP problem asks why we have never observed any CP-violating strong interactions. There are two possible answers to this question: either there is no observable strong CP violation or we just haven’t seen it. This section shows

¹In the terms of Section 2, discussions like those of Healey (2007, 2010) and Bain (2019) characterize the strong CP problem as asking why the Yang–Mills vacuum parameter θ_{YM} vanishes, when in fact the strong CP problem concerns the difference $\bar{\theta} = \theta_{\text{YM}} - \theta_{\text{Q}}$ of the Yang–Mills vacuum parameter and a phase on the quark mass matrix (Section 2.2). As such, one cannot solve the strong CP problem by showing that $\theta_{\text{YM}} = 0$, as Healey suggests.

This mischaracterization isn’t confined to the philosophical literature; physicists often elide θ_{YM} and $\bar{\theta}$ as well. For example, Healey cites a textbook by Rubakov, which also poses the strong CP problem in terms of θ_{YM} —at least in the main text (2002, 277). Rubakov clarifies in the appendix (2002, 417) that the problem in fact concerns $\bar{\theta}$.

that the second option can be ruled out in QCD. There are two mechanisms of CP violation in QCD: one in the Yang–Mills sector and one in the quark sector. The first is parametrized by a constant θ_{YM} and the second by a constant θ_{Q} . These can cancel: the Lagrangian as a whole is CP-invariant if and only if $\bar{\theta} = \theta_{\text{YM}} - \theta_{\text{Q}}$ vanishes. This parameter $\bar{\theta}$ can be brought into contact with experiment using an effective low-energy theory of strong interactions according to which $\bar{\theta}$ is proportional to the CP-violating electric dipole moment of the neutron. Observation places the neutron electric dipole moment very close to zero, so $\bar{\theta}$ is very close to zero as well—that is, QCD either respects CP or very nearly does. But we’re left with a new puzzle: why do θ_{YM} and θ_{Q} conspire to make $\bar{\theta}$ vanish?

2.1 The QCD Lagrangian

The strong sector of the Standard Model consists of an $SU(3)$ Yang–Mills field modeling the strong force and six Dirac fermions modeling the six quarks. For the sake of clarity and later discussion we’ll consider the more general case of an $SU(N_c)$ Yang–Mills theory and N_f Dirac fermions; call N_c the number of colors and N_f the number of flavors.

Our conventions for the $SU(N_c)$ gauge theory are as follows. The setting is 4-dimensional Minkowski spacetime \mathbb{M}^4 . Fix, once and for all, some natural coordinates on \mathbb{M}^4 . A configuration of the Yang–Mills theory can be specified by a vector potential, a one-form $A_\mu dx^\mu$ where A_μ takes values in $N_c \times N_c$ traceless hermitian matrices. Using this vector potential we define the gauge covariant derivative

$$D_\mu = \partial_\mu + i \frac{g}{\sqrt{N_c}} A_\mu$$

where ∂_μ is the standard flat derivative on \mathbb{M}^4 and g is a coupling constant. The curvature of this derivative, representing the field strength, is

$$F_{\mu\nu} = \partial_\mu A_\nu - \partial_\nu A_\mu + i \frac{g}{\sqrt{N_c}} [A_\mu, A_\nu]$$

If A_μ and A'_μ are two vector potentials then a gauge transformation sending the former to the latter is a smooth function $h : \mathbb{M}^4 \rightarrow SU(N_c)$ such that

$$A'_\mu = h A_\mu h^{-1} - i \frac{\sqrt{N_c}}{g} h \partial_\mu h^{-1}$$

When N_c is fixed its appearance in these expressions is a matter of convention; it can be eliminated by rescaling g or A_μ . Its appearance in the gauge potential and the field strength gives meaning to the expansion in terms of N_c^{-1} that we will consider later.

The QCD Lagrangian has four terms:

$$\mathcal{L}_{\text{QCD}} = -\frac{1}{2} \text{tr}(F_{\mu\nu} F^{\mu\nu}) + \frac{g^2 \theta_{\text{YM}}}{16\pi^2 N_c} \text{tr}(F_{\mu\nu} \tilde{F}^{\mu\nu}) + \bar{q} i \not{D} q - \bar{q} M e^{-i\theta_{\text{Q}} \gamma^5 / N_f} q$$

The first two terms describe the Yang–Mills sector and the other two the quark sector. The first term is the usual kinetic term for the Yang–Mills field. The second term, characterizing the Yang–Mills vacuum, depends on the Hodge dual $\tilde{F}^{\mu\nu} = \frac{1}{2} \epsilon^{\mu\nu\alpha\beta} F_{\alpha\beta}$ of the field strength tensor and a parameter θ_{YM} that

we call the “Yang–Mills vacuum parameter”. The quark sector involves an N_f -component vector q of quark fields charged under the Yang–Mills field. Each quark field is a color N_c -tuple, and each entry in this N_c -tuple is a Dirac spinor. The term $\bar{q}i\not{D}q$ is a gauge-invariant combination of the kinetic term for the quark fields and the interaction term between the quarks and Yang–Mills field, where the Dirac operator $\not{D} = \gamma^\mu D_\mu$ is a sum of the kinetic and interaction terms and acts componentwise on q . The final term is a mass term for the quarks built from a real diagonal $N_f \times N_f$ matrix M and an angle θ_Q .

The first three terms of the QCD Lagrangian are taken directly from the Standard Model; the fourth is an effective mass term coming from the Higgs mechanism. The Standard Model Lagrangian does not have any quark mass terms. It cannot, because quark mass terms violate electroweak gauge invariance: the $SU(2)$ Yang–Mills field in the electroweak sector couples only to the left-handed component of each quark, but the mass term

$$\bar{q}Me^{-i\theta_Q\gamma^5/N_f}q = \bar{q}_RMe^{-i\theta_Q/N_f}q_L + \bar{q}_LMe^{i\theta_Q/N_f}q_R$$

mixes the vector q_L of left-handed components with the vector q_R of right-handed components. Instead of a mass term, quarks receive masses from Yukawa interactions with the Higgs, and the couplings for these interactions may be complex. By performing a field redefinition we can absorb the Higgs’ vacuum expectation value and bring the Yukawa terms into the form we have assumed above. Since the Yukawa couplings will generally be complex, the angle θ_Q will generally be nonzero.

The Yang–Mills vacuum term and the quark mass term both generally violate CP symmetry. By the CPT theorem, CP violation is equivalent to violation of time-reversal symmetry. Since the Hodge dual $\tilde{F}_{\mu\nu}$ depends on the volume form $\epsilon^{\mu\nu\alpha\beta}$ it will pick up a sign under time reversal. Time reversal also conjugates complex coefficients, so it will send θ_Q to $-\theta_Q$ in the quark mass term. So the Yang–Mills vacuum term violates CP for nonzero θ_{YM} and the quark mass term violates CP for nonzero θ_Q .

2.2 The strong CP parameter

The quark mass term violates CP when θ_Q is nonzero, but θ_Q isn’t a physical parameter. The quantum theory is given by a path integral over all classical field configurations, meaning that we are free to redefine the fields by changing the variables of integration. In particular, we can redefine the quark fields so as to absorb the phase on the mass matrix. This redefinition is accompanied by a shift in θ_{YM} , so it is only the difference $\bar{\theta} = \theta_{YM} - \theta_Q$ that’s invariant under field redefinitions. That is, only $\bar{\theta}$ can have any physical significance.

We can set θ_Q to zero and eliminate CP violation in the quark sector by redefining the quark fields so as to absorb the phase on the mass matrix. Indeed, we have to perform a redefinition of this kind when we put the Yukawa terms of the Standard Model in the form of a quark mass term. Field redefinitions are permitted because we integrate over all possible field configurations in the path integral; a field redefinition is just a change in the variable of integration. The quantum theory is given by the path integral

$$Z = \int \mathcal{D}A \mathcal{D}q \mathcal{D}\bar{q} \exp\left(i \int d^4x \mathcal{L}_{\text{QCD}}\right)$$

By redefining the variables of integration q and \bar{q} we can eliminate the phase $\exp(i\theta_Q\gamma^5/N_f)$ in the quark mass term, but this requires some care.

A change of variables under an integral picks up a Jacobian determinant. In a familiar example, changing from rectangular to polar or spherical coordinates means making the replacements

$$dx dy \mapsto r dr d\phi \qquad dx dy dz \mapsto r^2 \sin \theta dr d\theta d\phi$$

respectively. The factor of r in the first expression and the factor of $r^2 \sin \theta$ in the second are the Jacobian determinants of the transformation from rectangular to polar coordinates. When the new coordinates in some transformation are linear functions of the old, the Jacobian determinant is a constant. And the redefinitions that put the Yukawa term in the form of a mass term are linear. So if the path integral were an integral in the usual sense we could simply set θ_Q to zero by absorbing a phase into the quark fields and making no changes to the path integral measure.

Because the path integral generally isn't an integral in the usual sense, a change of variables involves more than an ordinary Jacobian determinant (Weinberg, 1995, §22.2). In particular, notation like “ $\mathcal{D}A$ ” is heuristic: there is no Lebesgue measure on the space of field configurations, and integrals against other natural measures on this space tend to diverge. Finite values are extracted from a path integral with the help of a regulator—a method for suppressing the path integral's dependence on the details of physics at very high energies—and the effects of this regulator on the Jacobian determinant must be accounted for. The most naive regulator assigns measure zero to field configurations with momentum modes above some cutoff Λ . But this method is too naive: it isn't gauge-invariant, and it leads to pathologies like a divergent self-energy for the photon (Peskin and Schroeder, 1995, §7.5). Theories featuring gauge fields are more often regulated with dimensional regularization, which analytically continues the path integral measure to configurations defined on $(4 - \epsilon)$ -dimensional Minkowski space and takes the limit of vanishing ϵ at the end of the computation. However, absorbing the phase into the quark fields transforms the path integral measure as

$$\mathcal{D}A \mathcal{D}q \mathcal{D}\bar{q} \mapsto \mathcal{D}A \mathcal{D}q \mathcal{D}\bar{q} \exp\left(-i\theta_Q \int d^4x \operatorname{tr}(\gamma^5) \delta(x-x)\right)$$

where the divergent $\delta(x-x)$ requires regulation. Regulating this integral with dimensional regularization requires extending γ^5 to spacetimes with arbitrary complex dimension, and this cannot be done consistently while maintaining all of its relevant properties ('t Hooft and Veltman, 1972, §5). A field redefinition involving γ^5 requires another choice of regulator.

If we use a gauge-invariant regulator in four spacetime dimensions we obtain a nontrivial Jacobian determinant when absorbing the phase on the quark mass matrix. Following Fujikawa (1979, 1980) we impose a Gaussian cutoff, weighting the quark fields with $\exp(-\mathcal{D}^2/R^2)$ for a parameter R and taking the large R limit at the end of the calculation. For eigenmodes of \mathcal{D} with frequency ω this weight is $e^{-\omega^2/R^2}$, so this factor suppresses high-frequency modes of the quark fields. And it does so in a gauge-invariant way, because the weight is a function of the connection induced by the vector potential. Regulating the Jacobian

determinant in this way, we find that the path integral measure changes as

$$\mathcal{D}A \mathcal{D}q \mathcal{D}\bar{q} \mapsto \mathcal{D}A \mathcal{D}q \mathcal{D}\bar{q} \exp\left(-i \frac{\theta_Q g^2}{16\pi^2 N_c} \int d^4x \operatorname{tr}(F_{\mu\nu} \tilde{F}^{\mu\nu})\right)$$

The appearance of $F_{\mu\nu}$ isn't too surprising, given that our regulator is a gauge-invariant function of the gauge potential. More surprising, perhaps, is the fact that the change in θ_Q is compensated by a change in the Yang–Mills vacuum parameter: the regulated Jacobian determinant can be absorbed into the Lagrangian by replacing θ_{YM} with $\theta_{\text{YM}} - \theta_Q$.

Any observable CP violation can only depend on the relative angle between the CP-violating Yang–Mills term and the phase on the quark mass matrix. That is, observable CP violation is only possible if the parameter

$$\bar{\theta} = \theta_{\text{YM}} - \theta_Q$$

is nonzero. A field redefinition is just a change of variables; it cannot change the value of the integral. We can always make the mass matrix real by shifting the Yang–Mills vacuum parameter. And we can always set the Yang–Mills vacuum parameter to zero by rotating the phase on the quark mass matrix. The parameter $\bar{\theta}$ is a measure of CP violation that's independent of these redefinitions, analogous to the Jarlskog invariant in the electroweak sector. If—by some cosmic accident—the parameter $\bar{\theta}$ is zero, then making the mass matrix real will also make the Yang–Mills vacuum term vanish. In this case the Lagrangian will be CP invariant and so will all of its predictions.

We can use this result to resolve what you might call the “electromagnetic CP problem”, which asks why the electromagnetic force appears to respect CP symmetry. Since CP isn't a symmetry of nature, we need some explanation for the fact that electromagnetism respects it. Electromagnetism is the $N_c = 1$ case of our discussion, so electromagnetic CP violation is only possible if $\bar{\theta} = \theta_{\text{YM}} - \theta_Q$ is nonzero, with θ_{YM} the Yang–Mills vacuum parameter for the electromagnetic force. By a field redefinition we can put all of the CP violation in the Yang–Mills vacuum term of the Lagrangian, and by Stokes' theorem the only CP violation in the weighted path integral measure comes from the phase

$$\exp\left(i \frac{g^2 \bar{\theta}}{16\pi^2} \int_{\mathbb{M}^4} d^4x F_{\mu\nu} \tilde{F}^{\mu\nu}\right) = \exp\left(i \frac{g^2 \bar{\theta}}{16\pi^2} \int_{S_\infty^3} d^3x \epsilon^{\mu\nu\alpha} A_\mu F_{\nu\alpha}\right)$$

where S_∞^3 is the three-sphere at infinity and $\epsilon^{\mu\nu\alpha}$ is a volume form on S_∞^3 . If we assume that the field strength dies off at infinity—as is customary in quantum field theory—then the integrand on the right vanishes and this term drops out of the weighted path integral measure. So there is no CP violation in electromagnetism: up to a field redefinition the only CP-violating term of the action is the Yang–Mills vacuum term, and it vanishes.²

²We can also always neglect the vacuum term in the electroweak theory. As in the electromagnetic and strong case we can use a field redefinition to eliminate the vacuum term associated with the $SU(2)$ Yang–Mills vacuum term and put it in the Yukawa terms. Rotations of the right-handed quark fields can then be used to change the phase on the Yukawa couplings without reintroducing a vacuum parameter, since the weak isospin field couples chirally to the quark fields.

2.3 Low-energy QCD

We have little analytic control over the low-energy theory of the strong force. However, we can use the QCD Lagrangian to develop an effective theory of some low-energy phenomena. Low-energy effective theories are often obtained by integrating out the high-energy modes in the Lagrangian, but for QCD this approach does not work. Instead, we extract a low-energy effective theory from approximate symmetries of the Lagrangian. In particular, we hypothesize a spontaneously broken chiral symmetry that mixes the up and down quark. By Goldstone’s theorem this leads to three light bosons which we take to be the pions.

Extracting a low-energy theory of the strong force from the Standard Model is difficult. In the Standard Model the strong force is modeled as an $SU(3)$ Yang–Mills field. Its coupling to the quarks is described by the QCD Lagrangian with $N_c = 3$ and $N_f = 6$. But every particle in the low-energy spectrum of the theory is uncharged under the Yang–Mills field. So the low-energy degrees of freedom in the theory must be bound states of quarks whose charges cancel. This suggests that we should construct a low-energy theory by integrating out the high-energy degrees of freedom in the QCD Lagrangian. The lightest particles that interact with the strong force are the pions, whose mass is about 135 MeV, so we might try to consider physics below 1 GeV.³ Unfortunately, the strong force is strongly coupled at low energies and its vacuum is complicated. Analytically extracting a low-energy theory this way is prohibitively difficult, so we must advert to other methods.

Rather than integrating out the high-energy modes of the QCD Lagrangian, we can construct an effective theory of its low-energy modes by appealing to symmetry considerations. The masses of the two lightest quarks, the up and the down, differ by only a few MeV. In the context of GeV scale phenomena this difference is negligible, so rotations that mix the up and the down are an approximate symmetry of the Lagrangian. The only parameter in the QCD Lagrangian that depends on flavor is the mass, so if two quarks have the same mass then they are indistinguishable. More precisely, because q is an N_f -component vector the group $U(N_f)$ acts on it in the defining representation. The kinetic term $\bar{q}i\not{D}q$ is always invariant under this action, since \not{D} acts in the same way on every component of q and $U(N_f)$ just scrambles these components. If two quarks have the same mass then the mass matrix M commutes with the action of the $U(2)$ subgroup of $U(N_f)$ that mixes those two quarks, giving the mass term—and thus the Lagrangian—a $U(2)$ symmetry. If three quarks have the same mass then M commutes with the $U(3)$ subgroup that mixes those three quarks, giving the Lagrangian a $U(3)$ symmetry. And so on.

Not only do the up and the down have approximately the same mass, they are both approximately massless—only a few MeV—and this gives the Lagrangian further symmetry. If two quarks are massless then *a fortiori* they have the same mass, so the Lagrangian has the $U(2)$ flavor-mixing symmetry of the previous paragraph. But more is true. The kinetic term $\bar{q}i\not{D}q$ is chirally symmetric, so if we drop the mass terms then the left- and right-handed components of each quark field decouple. Two massless quarks thus give the Lagrangian a $U(2) \times U(2)$ symmetry, since the left- and right-handed components of each flavor

³All numerical quantities in this paper are taken from the 2018 Review of Particle Physics (Tanabashi et al., 2018), unless otherwise specified.

can be mixed independently. The $U(2)$ flavor-mixing symmetry of the previous paragraph acts on the left- and right-handed components in the same way, making it the diagonal subgroup of $U(2) \times U(2)$ —i.e., the subgroup of pairs in which the two entries coincide. Again, this situation generalizes straightforwardly if we have more massless quarks. If there are N_f massless quarks then the Lagrangian has a $U(N_f) \times U(N_f)$ symmetry that mixes the left- and right-handed components of these quark fields independently with the $U(N_f)$ action of the previous paragraph sitting on the diagonal.

We suppose that the approximate chiral symmetry of the up and down quarks is spontaneously broken. This is suggested first by the fact that pions do not come in parity doublets. If the vacuum shared the approximate $U(2) \times U(2)$ symmetry of the QCD Lagrangian then we should expect each light hadron to have a partner with opposite parity, but we don't see this parity doubling in nature. Second, lattice QCD computations give the following nonzero chiral condensate

$$\frac{1}{2}\langle\bar{q}q\rangle = \frac{1}{2}\langle\bar{q}_Lq_R + \bar{q}_Rq_L\rangle \approx (-250 \text{ MeV})^3$$

(Cichy et al., 2013, Table 7). Because the product $\bar{q}q$ mixes left- and right-handed components of q , it is only preserved by the diagonal subgroup of $U(2) \times U(2)$. Third, and most convincingly, the hypothesis of spontaneous symmetry breaking leads to a successful effective theory of pions.

Spontaneously broken symmetries inform us about low-energy physics because they imply the existence of light particles. Goldstone's theorem says that for every spontaneously broken, continuous, global, approximate symmetry of a Lorentz-invariant quantum theory there is one light spinless particle, the Nambu-Goldstone boson corresponding to that broken symmetry. A “symmetry of the quantum theory” means a symmetry of the effective action. For example, consider a theory of a set of scalar fields ϕ^i with action $S[\phi^i]$. The generating function for this theory is

$$Z[J_i] = \int \mathcal{D}\phi^i \exp\left(iS[\phi^i] + \int d^4x \phi^i J_i\right)$$

and the effective action is its Legendre transform

$$\Gamma[\phi^i] = -i \log Z[J_i] - \int d^4x \phi^i J_i$$

If $\Gamma[\phi^i]$ is invariant under a transformation generated by a matrix T_j^i then we can take two functional derivatives and evaluate at the minimum $\langle\phi^i\rangle$ of the effective action to give

$$\phi^i \mapsto \phi^i + i\epsilon T_j^i \phi^j \quad \frac{\delta^2\Gamma}{\delta\phi^k \delta\phi^i} T_j^i \langle\phi^j\rangle = 0$$

The second functional derivative of the effective action evaluated at its minimum is the two-point correlation function evaluated with vanishing external momenta and with all 1PI contributions included. In other words, it is the effective mass matrix. So this equation says that there is one mass eigenvector with eigenvalue zero for each nonvanishing $\langle\phi^i\rangle$. That is, for each spontaneously broken symmetry there is a massless boson. If the symmetry is only approximate then the right

Meson	π^0	π^\pm	K^\pm	K^0/\bar{K}^0	η	η'
Mass (MeV)	135	140	494	498	575	958

Table 1: The nine lightest pseudoscalar mesons

hand side of this equation only approximately vanishes, giving one light boson for each spontaneously broken generator (Weinberg, 1995, §19.3).

The chiral condensate breaks the up–down chiral flavor symmetry of the Lagrangian to a non-chiral flavor symmetry, giving three Nambu–Goldstone bosons. The $U(2) \times U(2)$ symmetry group of the QCD Lagrangian does not extend to a symmetry group of the effective action. Chiral rotations

$$q \mapsto q + i\epsilon\gamma^5 q$$

are a symmetry of the Lagrangian, but as we saw in Section 2.2 they aren’t a symmetry of the path integral measure, so they aren’t a symmetry of the effective action either. So Goldstone’s theorem doesn’t apply to the $U(2) \times U(2)$ symmetry of the QCD Lagrangian. However, it does apply if we restrict attention to the $SU(2) \times SU(2)$ subgroup, since chiral rotations lie outside this subgroup. The Lie algebra of this group is six dimensional, and the Lie algebra of the unbroken diagonal subgroup has three dimensions, so Goldstone’s theorem gives three Nambu–Goldstone bosons. The broken generators are odd under parity, so the Nambu–Goldstone bosons must be pseudoscalars; the nine lightest candidates are listed in Table 1. The most important feature of this table is the striation of the particle tuplet masses: there is a triplet of pions with masses of about 140 MeV, a quadruplet of kaons with masses near 500 MeV, the η whose mass is slightly larger than the kaons, and the η' whose mass is nearly 1 GeV. A theory of low-energy QCD should account for this mass spectrum in terms of parameters appearing in the QCD Lagrangian. We will return to this in Section 4.2.

We suppose that the three pions are the three Nambu–Goldstone bosons of the spontaneously broken $SU(2) \times SU(2)$ symmetry. The cash value of this supposition is that low-energy pion physics can be effectively described using the most general theory with the symmetry-breaking pattern of the QCD Lagrangian (Scherer and Schindler, 2012). The standard construction gives a theory of an $SU(2)$ -valued field

$$U_2 = \exp\left(\frac{i}{F_\pi} \begin{pmatrix} \pi^0 & \sqrt{2}\pi^- \\ \sqrt{2}\pi^+ & -\pi^0 \end{pmatrix}\right)$$

where the pion decay constant F_π has dimensions of mass.⁴ To lowest order in U_2 , the effective Lagrangian is the chiral Lagrangian

$$\mathcal{L}_{\chi,2} = \frac{1}{4}F_\pi^2 \text{tr}(\partial^\mu U_2^\dagger \partial_\mu U_2) + \frac{1}{2}B_\pi F_\pi^2 \text{tr}\left(e^{-i\bar{\theta}/2} M U_2 + e^{i\bar{\theta}/2} M U_2^\dagger\right)$$

The first of these terms is the lowest-order term permitted by the symmetries of the system, and the numerical value $F_\pi = 92.07$ MeV is determined by measuring

⁴By the “standard construction” I mean that of Coleman et al. (1969) and Callan Jr. et al. (1969); see also Weinberg (1995, §19.6).

rates of leptonic pion decay. The second term explicitly breaks the $SU(2) \times SU(2)$ symmetry while leaving its diagonal subgroup intact. The constant B_π can be interpreted by equating the vacuum energies of the QCD and chiral Lagrangians. To zeroth order in $\bar{\theta}$ this gives

$$\frac{1}{2}\langle\bar{q}q\rangle = -B_\pi F_\pi^2$$

So B_π is proportional to the chiral condensate if $\bar{\theta}$ is small.

This treatment of pions as Nambu–Goldstone bosons has been extremely successful. A somewhat simple success of the model is an explanation of the pion masses in Table 1. Expanding the chiral Lagrangian in the pion fields about the minimum of U_2 gives the mass

$$m_\pi^2 = B_\pi \sqrt{m_u^2 + m_d^2 + 2m_u m_d \cos \bar{\theta}}$$

The lightness of the pions is therefore a result of the lightness of the up and down quarks. The pions exhaust the bound quark–antiquark states of the up and down quarks, so any other pseudoscalar mesons must include some heavier quarks. And indeed, the pions are significantly lighter than the next lightest pseudoscalar mesons. They are also close to one another, and the remaining discrepancy is mostly accounted for by electromagnetic contributions. The reproduction of this result, the Gell-Mann–Oakes–Renner relation, is one success of the Nambu–Goldstone model. Others include various aspects of soft pion physics, including leptonic pion decays, pion photoproduction, and pion–pion scattering.⁵ We can also use it to extract an observable consequence of $\bar{\theta}$.

2.4 The neutron electric dipole moment

The neutron can acquire an electric dipole moment (EDM) through its interactions with the pion. Because a neutron EDM would violate CP, any interactions responsible for it must come from the CP-violating terms in the QCD Lagrangian, making the neutron EDM a measure of $\bar{\theta}$. It is a good measure: it is tightly constrained by observation and is proportional to $\bar{\theta}$, thereby tightly constraining the latter. The neutron has no observable electric dipole moment, so $\bar{\theta}$ must be vanishingly small. This is surprising: why would the strong force and the quark sector conspire to make $\bar{\theta}$ vanish? This is the usual formulation of the strong CP problem.

A particle with a permanent EDM violates CP. Dipole couplings of a Dirac spinor ψ to the electromagnetic field are described by terms of the form

$$\mathcal{L}_{\text{dipole}} = \mu_\psi F_{\mu\nu} \bar{\psi} S^{\mu\nu} \psi + d_\psi \tilde{F}_{\mu\nu} \bar{\psi} S^{\mu\nu} \psi$$

with $F_{\mu\nu}$ the electromagnetic field strength tensor and $S^{\mu\nu}$ the generators of the bispinor representation of the Lorentz group. The couplings μ_ψ and d_ψ are the magnetic and electric dipole moments of ψ , respectively. The first term is invariant under the entire Lorentz group, but the $\tilde{F}_{\mu\nu}$ in the electric term violates CP symmetry just as it does in the QCD Lagrangian.

⁵See Donoghue et al. (1992, Ch. VI), Scherer and Schindler (2012, Ch. 3), and Weinberg (1995, §19.4) for further discussion of this model.

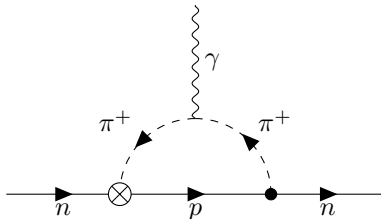


Figure 1: A CP-violating contribution to the neutron electric dipole moment

If $\bar{\theta}$ is nonzero then interactions with charged pions give the neutron an effective EDM (Srednicki, 2007, §94). To compute this EDM we extend the chiral Lagrangian of Section 2.3 to include the proton and neutron in a nucleon doublet N that transforms in the defining representation of $SU(2)$. Adding the leading-order terms involving N and U_2 and expanding in the pion fields and $\bar{\theta}$, we find that the leading interaction terms between the pions and nucleons are of the form

$$\mathcal{L}_{\pi N} = \bar{N}(g_{\pi N}i\gamma^5 + \bar{g}_{\pi N}\bar{\theta})\pi N$$

with π the $\mathfrak{su}(2)$ -valued field that generates U_2 . The constants $g_{\pi N}$ and $\bar{g}_{\pi N}$ can be fixed using pion–neutron scattering cross sections and baryon mass differences, respectively (Crewther et al., 1979). Because the pions are pseudoscalars the first term of this Lagrangian is CP-invariant and the second is not. So, as expected, CP violation is parametrized by $\bar{\theta}$. This Lagrangian leads to CP-violating processes like the one in Fig. 1, which involves one CP-violating vertex (marked by a cross) and one CP-preserving vertex. Making the first vertex CP-preserving and the second CP-violating gives another contribution to the neutron EDM.

The neutron’s EDM is obtained by matching matrix elements. In the soft photon limit the electric term in $\mathcal{L}_{\text{dipole}}$ gives a matrix element of the form

$$iM = 2d_n \varepsilon_\mu^*(q) \bar{u}(p') S^{\mu\nu} q_\nu i\gamma^5 u(p)$$

with p and p' the incoming and outgoing momenta of the neutron, respectively, and q the momentum of the incoming photon. On the other hand, the diagram in Fig. 1 and its vertex-swapped partner give the matrix element

$$iM = |\bar{\theta}| \frac{eg_{\pi N}\bar{g}_{\pi N}}{4\pi^2 m_N} \log\left(\frac{\Lambda^2}{m_\pi^2}\right) \varepsilon_\mu^*(q) \bar{u}(p') S^{\mu\nu} q_\nu i\gamma^5 u(p)$$

with m_N the average mass of the nucleons and Λ a momentum cutoff for the loop. Matching the coefficients on these terms and taking the cutoff Λ to be $4\pi F_\pi$, comparison with experiment gives

$$|\bar{\theta}| < 7.6 \times 10^{-11}$$

More conservative estimates place $|\bar{\theta}|$ as high as 10^{-9} (Vicari and Panagopoulos, 2009, §7.1).

2.5 The strong CP problem

The vanishingly small value of $\bar{\theta}$ amounts to outrageously delicate cancellation between the strong force and quark sectors. This is the usual guise of the strong

CP problem: why is $\bar{\theta}$ so small? Healey’s criticism of the Standard Model doesn’t turn on precisely what the problem here is, so we can mostly remain at this level of generality. For the sake of completeness I briefly sketch different ways you might make the problem more precise and describe the main contenders in the search for a solution.

Let me recapitulate the problem, now that all the moving parts are laid out. According to the Standard Model, strong interactions are described by a Lagrangian of the form

$$\mathcal{L}_{\text{QCD}} = \mathcal{L}_{\text{CP}} + \frac{g^2 \theta_{\text{YM}}}{16\pi^2 N_c} \text{tr}(F_{\mu\nu} \tilde{F}^{\mu\nu}) - \bar{q} M e^{-i\theta_{\text{Q}} \gamma^5 / N_f} q$$

with \mathcal{L}_{CP} consisting of CP-preserving terms. Assuming a gauge-invariant regulator, neither θ_{YM} nor θ_{Q} has any coordinate-invariant meaning, but their difference $\bar{\theta}$ does. The quantum theory will be CP invariant just in case $\bar{\theta}$ vanishes. At low energies QCD is strongly coupled, so we cannot directly extract predictions of CP violation from a nonzero $\bar{\theta}$. However, we can use the symmetry breaking pattern of the QCD Lagrangian to construct an effective low-energy theory of pions and nucleons, and this theory predicts that the neutron EDM is proportional to $\bar{\theta}$. Measurements of this EDM show that $\bar{\theta}$ vanishes to one part in ten billion. So the strong force is CP symmetric—or at least very nearly so.

Prima facie, this situation is surprising, and perhaps a sign that new physics is required. The parameter $\bar{\theta}$ can take any value between 0 and 2π , since θ_{Q} is a complex argument. And there doesn’t seem to be any reason for it to take one value rather than another. So why should it happen to take precisely the value that gives strong CP symmetry? This is a question of conspiracy, not numerology. The question isn’t why some fundamental parameter takes on an interesting value, in the way that the fine structure constant is $1/137$ or the proton–electron mass ratio is $6\pi^5$. The question is why two constants θ_{YM} and θ_{Q} from apparently unrelated terms in the Lagrangian happen to take on the same value, given that this is necessary and sufficient for CP violation in the strong sector.

Filling in the explanatory context a bit more, we might ask why $\bar{\theta}$ vanishes given that CP isn’t a symmetry of nature. We have seen CP violation in weak decays involving kaons, B mesons, and D mesons. In the Standard Model this CP violation is accounted for by the Kobayashi–Maskawa mechanism: with three generations of quarks and complex Yukawa couplings there can be CP-violating flavor mixing in charged weak interactions. The possibility of complex Yukawa couplings is required for this mechanism. And this means there’s no reason for θ_{Q} to vanish when we put the Yukawa terms in the form of quark mass terms. But if θ_{Q} doesn’t vanish, what reason could there be for θ_{YM} to take on the same non-vanishing value? It’s possible to have Yukawa couplings that permit weak CP violation while θ_{Q} vanishes, but this is just shifting the explanatory problem from $\bar{\theta}$ to the Yukawa couplings. And it still gives no reason for θ_{YM} to take on the same value as θ_{Q} .

The smallness of $\bar{\theta}$ is often included in lists of “naturalness” problems with the Standard Model, which have been of interest in recent philosophical literature. It certainly is one on some understandings of naturalness problems. For example, in a theory that’s natural in ’t Hooft’s (1979) sense a parameter like $\bar{\theta}$ can only be small if setting it to zero would increase the symmetry of the theory. Since CP isn’t a symmetry of the Standard Model the parameter $\bar{\theta}$ is unnatural in

this sense. But every unhappy theory is unhappy in its own way, and the strong CP problem isn't obviously problematic in a way that matters. Williams (2015) argues that naturalness problems are problems because they violate a “central dogma” of effective field theory, “a prohibition on sensitive correlations between widely separated physical scales” (2015, 95). And the parameter $\bar{\theta}$ doesn't seem to violate this prohibition. In the first place, the value of $\bar{\theta}$ doesn't make much difference to large-scale physics: the electric dipole moments it induces will slightly shift atomic energy levels, but its effects will be so shielded by the lightness of the quarks that these shifts will make no qualitative difference to, for example, the stability of any elements. Nor does the value of $\bar{\theta}$ depend sensitively on the high-energy cutoff of the theory. The parameter θ_{YM} does not renormalize, since the vacuum Yang–Mills term vanishes in perturbation theory. The parameter θ_{Q} does receive corrections, but they do not occur before the seven loop level (Ellis and Gaillard, 1979).

The search for a solution to this problem is going strong. The most popular explanation of $\bar{\theta}$'s value, due to Peccei and Quinn (1977a,b), posits a new spontaneously broken $U(1)$ symmetry whose Nambu–Goldstone boson, called the axion, couples to the Yang–Mills vacuum term. After rotating $\bar{\theta}$ into the vacuum term, the axion acquires a vacuum expectation value that cancels $\bar{\theta}$ out. Axion searches are ever-increasing, in part because the axion is a dark matter candidate. However, there have been no signs so far. A second solution proposes that CP is only spontaneously broken, not explicitly broken (Barr, 1984; Nelson, 1984). This would protect $\bar{\theta}$ from being too large, but a solution along these lines must explain how the spontaneous symmetry breaking gives a phase on the Yukawa couplings that's large enough to account for weak CP violation. Such solutions generally involve fine-tuning or increasingly speculative physics beyond the Standard Model.

3 Large gauge transformations

If gauge transformations are eliminable redundancies then the reasoning of Section 2 is mistaken. The Yang–Mills vacuum term is not preserved by all gauge transformations. If the eliminative view of gauge transformations is right, this means that the Yang–Mills vacuum term is physically meaningless. If gauge transformations are redundancies then mathematical differences between gauge-equivalent configurations can't reflect physical differences. So the value of the Yang–Mills vacuum term can't represent any physical fact. Healey (2007, 2010) argues that the strong CP problem dissolves when we recognize this. If we treat the Yang–Mills vacuum term in accordance with the eliminative view then θ_{YM} is either ill-defined or makes no difference to the physics. Either way, some step in Section 2 misfires, and the strong CP problem never arises.

3.1 Giving global sense to the classical action

The site of disagreement is the classical action, which appears in the weight of the path integral measure. For simplicity, consider the pure Yang–Mills sector, where the Lagrangian is

$$\mathcal{L}_{\text{YM}} = -\frac{1}{2} \text{tr}(F_{\mu\nu} F^{\mu\nu}) + \frac{g^2 \theta_{\text{YM}}}{16\pi^2 N_c} \text{tr}(F_{\mu\nu} \tilde{F}^{\mu\nu})$$

Formally, the action of some Yang–Mills configuration is given by the expression

$$S_{\text{YM}} = \int_{\mathbb{M}^4} d^4x \mathcal{L}_{\text{YM}}$$

considered as a real-valued function on the space of vector potentials. However, this formal expression is ill-defined: since \mathbb{M}^4 isn't compact, this integral diverges for most vector potentials. Broadly speaking, there are two ways to reckon with this failure: either generalize the framework to include variational problems for which the action is only locally defined or restrict the configuration space to configurations for which the action is well-defined. In either case we find that the Yang–Mills vacuum term varies under some gauge transformations.

When it comes to the classical equation of motion we can forget about the variational framework and use the usual coordinate expression for the Euler–Lagrange equation to write down an equation of motion formally related to \mathcal{L}_{YM} . The resulting Yang–Mills equation of motion is

$$D_\mu F^{\mu\nu} = \partial_\mu F^{\mu\lambda} + i \frac{g}{\sqrt{N_c}} [A_\mu, F_{\mu\nu}] = 0$$

Over a compact spacetime manifold a configuration A_μ satisfies this equation just in case it's a stationary point of the Yang–Mills action. In that context the Yang–Mills equation really is the Euler–Lagrange equation of the Yang–Mills Lagrangian density, and the usual features of the variational framework apply—for example, Noether's theorems give conserved currents and charges. If we drop the variational framework then these features are absent, and the Yang–Mills equation has no mathematical relation to the Yang–Mills Lagrangian. Since the Lagrangian plays no role in this theory we might as well forget about it. In particular, since the vacuum parameter θ_{YM} doesn't appear in the Yang–Mills equation of motion it has no dynamical significance if we restrict our attention in this way.

We can regain a connection between the Yang–Mills Lagrangian and the equation of motion if we consider local variations. The Yang–Mills equation is a differential equation expressing a condition at each point of the spacetime manifold. Whether it holds at any given point is a fact about the infinitesimal neighborhood of that point. So a configuration satisfies the Yang–Mills equation everywhere on \mathbb{M}^4 just in case its restriction to any compact subregion K of \mathbb{M}^4 is a stationary point of the local action

$$S_{\text{YM}}^K = \int_K d^4x \mathcal{L}_{\text{YM}}$$

Using the usual differential-geometric methods, the variational framework can be extended to any smooth manifold from its compact submanifolds.⁶ This recovers a sense in which the Yang–Mills equation is the Euler–Lagrange equation of the Yang–Mills Lagrangian density. It also allows for generalizations of, for example, Noether's theorems.

In the local variational framework the vacuum parameter θ_{YM} makes a difference only to the action and not to the equation of motion. For any compact region K the value of S_{YM}^K depends on θ_{YM} , and changing the value of θ_{YM} will

⁶See Bleecker (1981), Saunders (1989), and Anderson (1992) for increasingly mathematical versions of this generalization.

change the action of a configuration within K . However, such changes make no difference to the equation of motion because the vacuum term is invisible to all variations. If we define

$$C_{\mu\nu\alpha} = A_\mu F_{\nu\alpha} - i \frac{2g}{3\sqrt{N_c}} A_\mu A_\nu A_\alpha$$

then by Stokes' theorem we have

$$\int_K d^4x \operatorname{tr}(F_{\mu\nu} \tilde{F}^{\mu\nu}) = \int_{\partial K} d^3x \epsilon^{\mu\nu\alpha} \operatorname{tr}(C_{\mu\nu\alpha})$$

where $\epsilon^{\mu\nu\alpha}$ is the volume form on the boundary ∂K of K . Any variation over K must vanish on the boundary, so the integral on the right—and therefore the vacuum term in the Yang–Mills action—will be left unchanged by any variation. The Euler–Lagrange equation of this action is thus independent of the vacuum term, hence insensitive to θ_{YM} .

The vacuum term appearing in the local action varies under some transformations. Under a gauge transformation h the current $C_{\mu\nu\alpha}$ is sent to the current

$$C_{\mu\nu\alpha} + 2i \frac{\sqrt{N_c}}{g} \partial_\mu (A_\nu h^{-1} \partial_\alpha h) + \frac{2N_c}{3g^2} (h \partial_\mu h^{-1}) (h \partial_\nu h^{-1}) (h \partial_\alpha h^{-1})$$

Integrating this over ∂K kills the second term by Stokes' theorem, so applying the gauge transformation h sends the action over K to

$$S_{\text{YM}}^K + \frac{\theta_{\text{YM}}}{24\pi^2} \int_{\partial K} d^3x \epsilon^{\mu\nu\alpha} \operatorname{tr}((h \partial_\mu h^{-1}) (h \partial_\nu h^{-1}) (h \partial_\alpha h^{-1}))$$

Whether the action is invariant under h —i.e., whether the second term vanishes—depends on what exactly we take h to be. Since the domain of integration is ∂K it's natural to consider gauge transformations corresponding to smooth functions of the form $h : \partial K \rightarrow SU(N_c)$. If h is a gauge transformation of this kind then it will generally change the action; the new term does not vanish. However, if h is a smooth function defined on all of Minkowski space then the new term vanishes and S_{YM}^K is unchanged. In other words, the locally defined action is invariant under gauge transformations that smoothly extend to all of \mathbb{M}^4 but may change under locally defined gauge transformations.

The disagreement over the size distinction is a disagreement about what to say in this situation. Strictly speaking there is no violation of gauge invariance: if A_μ and A'_μ are gauge equivalent configurations over \mathbb{M}^4 then they have the same action over every compact region. Nevertheless, you might think that a gauge transformation defined on the boundary ∂K or in some neighborhood of it should be interpreted just the same as any other gauge transformation. As Healey puts it, “local’ gauge symmetry is a purely formal feature of a theory,” and it follows from this that “a gauge transformation cannot connect representations of physically distinct situations,” even if it is only defined on the boundary (2007, 176). If you accept this principle then the vacuum term in the Yang–Mills Lagrangian is ill-defined and should be excluded from the Lagrangian. This is a rejection of the size distinction. Alternatively, you could accept the size distinction and conclude that gauge transformations defined on

∂K are symmetries relating distinct configurations with many—but not all—of the same properties.

Taking a different approach to the calculus of variations on \mathbb{M}^4 doesn't avoid the problem. Rather than adverting to locally defined actions, we can ensure that the action is globally defined by imposing some conditions on the configuration that make the integral converge (Marsden and Hughes, 1994, §5.1). Many different conditions could do this job; here the disagreement over the size distinction is a disagreement about which of these conditions are legitimate. The most naive condition would be to demand that A_μ be compactly supported or that it vanish sufficiently quickly at infinity. However, a moment's reflection shows that this condition is physically meaningless. Any compactly supported configuration is gauge equivalent over \mathbb{M}^4 to one that isn't—just apply a gauge transformation $h : \mathbb{M}^4 \rightarrow SU(N_c)$ such that $h\partial_\mu h^{-1}$ isn't compactly supported—and so no physically significant property can distinguish between configurations that are compactly supported and those that aren't. The integral of the Lagrangian density over all of \mathbb{M}^4 may give a well-defined action on the space of compactly supported configurations, but there's no reason to be interested in this space.

A full reckoning with the problem of gauge-invariant boundary conditions deserves a long treatment; I pursue it elsewhere. For the purposes of the present discussion, consider the space of configurations that are pure gauge outside of some compact region. That is, a vector potential A_μ represents a configuration in our restricted configuration space if there is some compact region K of \mathbb{M}^4 and some smooth function $h : \bar{K} \rightarrow SU(N_c)$ on the complement $\bar{K} = \mathbb{M}^4 \setminus K$ of K such that

$$0 = hA_\mu h^{-1} - i \frac{\sqrt{N_c}}{g} h\partial_\mu h^{-1}$$

on \bar{K} . This condition does not distinguish between gauge equivalent vector potentials because gauge equivalence is transitive. And it is plausibly a configuration space of some physical interest: if a configuration is pure gauge then it has vanishing field strength, hence exerts no forces on test particles.⁷

The action is globally defined on the space of configurations with compact gauge-invariant support, but the vacuum term again leads to a violation of local gauge invariance. If A_μ is pure gauge outside of some compact region K then its global action coincides with the action assigned to K in the local action formalism, since $F_{\mu\nu}$ will vanish outside of K . Just as in the local action formalism, the second term of this action is an integral over a total derivative. And just as before, a gauge transformation $h : \partial K \rightarrow SU(N_c)$ will generally change its value. This means, in particular, that the value of the vacuum term can be nonzero. This is true even though the second term can be written as an integral over ∂K , where the configuration is gauge equivalent to zero by hypothesis.

The vacuum term in the Yang–Mills Lagrangian can only contribute to the action in a well-defined way if some gauge transformations relate distinct states of affairs. The disagreement between the eliminative view and the Standard

⁷The appropriate notion of gauge transformation $(A_\mu, K, h) \rightarrow (A'_\mu, K', h')$ is a gauge transformation $k : \mathbb{M}^4 \rightarrow SU(N_c)$ from A_μ to A'_μ such that $h' \cdot k = h$ on $K \cap K'$. This choice gives the homotopy fiber over 0 of the map sending each vector potential to its boundary condition.

Model is therefore a disagreement over how to evaluate the integral

$$\int d^4x \operatorname{tr}(F_{\mu\nu}\tilde{F}^{\mu\nu})$$

in the classical action. Accepting the size distinction means allowing this integral to be nonzero. If we evaluate this integral over some compact region K —either because we are considering a locally defined action or because we have imposed boundary conditions—then a gauge transformation on the boundary ∂K will change the value of this integral. This gauge transformation is “large” in the sense that it is nontrivial on the boundary of the region of integration. In particular, if we demand that the configuration be pure gauge at infinity then a large gauge transformation is one that is nontrivial at infinity, recovering the usual statement of the size distinction. If we reject the size distinction and demand that gauge transformations on the boundary be treated just as gauge transformations elsewhere then this integral is ill-defined. The vacuum Yang–Mills term must therefore be excluded.

3.2 Holonomy theories

The eliminative view of gauge transformations is incompatible with size distinction, hence incompatible with the reasoning in Section 2. This is a reason to reject the eliminative view as an interpretation of our best theory of high-energy physics: it doesn’t describe that theory. But it might be that, all things considered, we should reject the size distinction and revise the Standard Model. This is Healey’s position. His objection to the size distinction is one part of a much larger argument for an interpretation of Yang–Mills theory in terms of properties attaching to curves in spacetime. On this interpretation the size distinction is theoretically unjustified—it can’t even be expressed—and the strong CP problem harmlessly dissolves. In what follows I will mostly set aside the details of Healey’s positive view, however. The conflict between the eliminative view and the Standard Model does not depend on the details of a particular eliminative view. And, as I argue in this section, continuing to talk in field-theoretic rather than curve-theoretic terms does not prejudge the issue.

On Healey’s interpretation of Yang–Mills theory the fundamental quantities of the theory are attached to loops in spacetime. In the case of electromagnetism, for example, there is a complex phase attached to each loop—the holonomy of the principal $U(1)$ -connection that represents the field. This view is inspired in part by the work of Wu and Yang (1975). As they say, there’s a sense in which the holonomy “provides a complete description that is neither too much nor too little” (1975, 3846). The electromagnetic field strength tensor $F_{\mu\nu}$ contains too little information, as demonstrated by the Aharonov–Bohm experiment. An electron moving outside of a solenoid can detect magnetic flux through the solenoid, even when $F_{\mu\nu}$ vanishes in the region through which the electron moves. So the electron must be detecting some facts that aren’t encoded in the field strength tensor, making the description provided by this tensor “too little”. The electromagnetic vector potential A_μ contains enough information to account for the electron’s behavior, but it’s “too much”: different vector potentials related by a gauge transformation have the same observable features. The holonomy an

electromagnetic potential A_μ assigns to a loop γ , given by

$$\exp\left(i \int_\gamma A_\mu dx^\mu\right)$$

is invariant under gauge transformations and provides enough information in the exterior region of the solenoid to distinguish between different values of the magnetic flux through the solenoid. So, at least when it comes to the Aharonov–Bohm effect, the holonomy is just the right amount of information. Healey’s interpretation of Yang–Mills theory develops and extends this idea to all Yang–Mills theories.

The discussion of Section 3.1 formulated Yang–Mills theory in terms of vector potentials, but it’s compatible with Healey’s view. If we take a configuration of Yang–Mills theory to be a smooth map assigning group elements to loops in spacetime then it’s still true that a configuration can be specified by a vector potential A_μ ; this potential gives a map H that acts on loops as

$$H(\gamma) = \mathcal{P}\exp\left(i \frac{g}{\sqrt{N_c}} \int_\gamma A_\mu dx^\mu\right)$$

where $\mathcal{P}\exp$ is the path-ordered exponential. The vector potential A_μ can therefore be understood as a choice of representative for the holonomy map it induces—a choice that’s eliminable, at least in principle. The Lagrangian involves only $F_{\mu\nu}$, not A_μ , and we can define $F_{\mu\nu}$ directly from a holonomy map. Given any such map H , the value of $F_{\mu\nu}$ at some spacetime point x is defined by

$$H(\gamma) = 1 - \epsilon^2 F_{\mu\nu}$$

where γ is an infinitesimal square in the x^μ – x^ν plane that’s based at x and has area ϵ^2 . The Lagrangian, and in particular the vacuum term, can therefore be defined directly in terms of the holonomy map.

Healey’s view is incompatible with the size distinction because the vacuum term is ill-defined. Large gauge transformations can change the value of the vacuum term, but Healey’s loop theory is—by design—invariant under all gauge transformations. So there is no way to consistently assign a value to the vacuum term in this theory. As Healey puts it, “there is no possibility of introducing a parameter $[\theta_{\text{YM}}]$ ” (Healey, 2007, 198).⁸ And this is his solution to the strong CP problem: we can’t consistently include the CP-violating term in the action, so there is no way for the strong force to violate CP, and this is why we haven’t seen any strong CP violation. This is a theoretical virtue of this loop-theoretic formulation of Yang–Mills theory: it solves a problem that its rivals cannot.

By formalizing the size distinction in terms of the Yang–Mills vacuum term I aim to isolate one feature of Healey’s view while holding the rest of the view fixed. I argue elsewhere that the size distinction is in fact orthogonal to the difference between field-theoretic and curve-theoretic formulations: some curve-theoretic formulations feature the size distinction and some don’t. In particular, the size distinction arises in the theory Wu and Yang actually study—after discussing the holonomy map as motivation, they set it aside and consider a theory in

⁸Strictly speaking, Healey is here referring to θ_{YM} as it appears in the algebraic approach to large gauge transformations. These θ_{YM} s are essentially equivalent, as he notes in his original discussion (2007, 197) and again, slightly more explicitly, in a later treatment (2010, 115).

which group elements are assigned to all curves in spacetime, not just closed loops. This formalism is fully equivalent to the vector potential formalism, hence includes the size distinction; Healey’s loop theory is not. Unfortunately, these two kinds of formulations are regularly lumped together. The loop theory is more consistent with Healey’s general gauge-elimination project—unlike the path theory, the loop theory lacks all talk of gauge and the size distinction—and so it is plausibly the one we should attribute to Healey. However, framing the dispute over the size distinction using the vacuum term means that we needn’t settle the issue here.

Formulating the size distinction as a dispute over the Yang–Mills vacuum term captures the physical consequences of the distinction while remaining agnostic about its status with respect to first principles. Once we have identified the vacuum term as the locus of disagreement its physical consequences can be extracted in the standard way, as we will see in the following sections. This is not the only way that the size distinction finds formal expression; for example, in canonical quantization it arises because Gauss’s law must be imposed on the Hilbert space as an operator equation. The path integral approach is just particularly convenient for our phenomena of interest. Phrasing it this way doesn’t beg the question against Healey: despite the use of vector potentials in setting up the Yang–Mills Lagrangian, the resulting expression can be interpreted in terms of a loop-based formulation of Yang–Mills theory. And, as he argues, his view implies that the vacuum term cannot contribute to the action. But its contributions are necessary if we want the right predictions.

4 The need for the vacuum term

Healey argues that the reasoning in Section 2 is mistaken because it is insufficiently gauge-invariant. Large gauge transformations can change the value of the Yang–Mills vacuum term, and this will change the weight on the path integral measure, so something has gone wrong. Healey offers two different diagnoses of the problem, and these suggest two different cures. On the first diagnosis we made a mistake by even allowing the vacuum Yang–Mills term in the first place. On Healey’s interpretation of gauge transformations, “there is no possibility of introducing a parameter” θ_{YM} in the QCD Lagrangian (Healey, 2007, 198). This suggests that the Yang–Mills vacuum term should simply be excluded from the Lagrangian. Alternatively, the problem might be due to the way that we evaluated the vacuum term: Healey claims that “once formulated, the loop representation will be equivalent to the usual connection representation with $[\theta_{\text{YM}}] = 0$ ” (Healey, 2007, 198). We might read this as the claim that the Yang–Mills vacuum term can be included, and the loop representation shows that the integral will always vanish, giving a theory equivalent to one in which θ_{YM} is zero. Neither of these proposals will work, however. We need the vacuum term to exist and have nonzero integral if we’re going to account for low-energy hadron physics.

The next two sections modify the theory of pions with each of Healey’s suggestions in turn. I take it that Healey would have no objection to the developments of Section 2 once the Yang–Mills vacuum term is excised from the QCD Lagrangian or vanishes in the QCD action and that the resulting theory is the one whose empirical content he endorses. Healey says in his

presentation of the strong CP problem that “[e]xperimental tests have shown that $|\bar{\theta}| \leq 10^{-10}$ ” (2007, 197), and as we saw in Section 2.4 this constraint is obtained by relating $\bar{\theta}$ to the neutron electric dipole moment using the pion theory of Section 2.3. To be sure, the derivation in Section 2.4 is only one way to compute $d_n/\bar{\theta}$; this ratio can also be derived using QCD sum rules, various bag models, current algebra, and more (Vicari and Panagopoulos, 2009, §7.1). Healey doesn’t cite a particular derivation of the ratio. But the results in this section can be derived in many ways, and they will apply to whatever theory of low-energy QCD Healey prefers. At any rate, Healey’s dissolution of the strong CP problem concerns the parameter θ_{YM} and not the theory of low-energy QCD more generally.

We have general grounds for thinking that Healey would accept the theory of pions in Section 2.3, but the role of spontaneous symmetry breaking in this theory deserves particular comment. Recall that a theory exhibits spontaneous symmetry breaking if there is some symmetry of the laws that is not a symmetry of the vacuum. In Section 2.3 we noted that in the limit of massless quarks the QCD Lagrangian has a $U(N_f) \times U(N_f)$ phase symmetry, where an element (g_L, g_R) of $U(N_f) \times U(N_f)$ acts as g_L on the left-handed components of the quark fields and g_R on the right-handed components, in both cases in the defining representation. We then supposed that the vacuum expectation value $\langle \bar{q}q \rangle$ is nonzero. An element (g_L, g_R) of $U(N_f) \times U(N_f)$ only preserves this vacuum expectation value if $g_L = g_R$ —that is, only if (g_L, g_R) belongs to the diagonal $U(N_f)$ subgroup of $U(N_f) \times U(N_f)$. So there is a $U(N_f) \times U(N_f)$ symmetry of the laws that isn’t a symmetry of the ground state, meaning that the symmetry is spontaneously broken. Whatever else is the case, then, the predictions of QCD will be the predictions of a theory with a spontaneously broken $U(N_f) \times U(N_f)$ symmetry. The theory of pions in Section 2.3 is just the most general theory with this feature that has the parameters of the QCD Lagrangian, so any predictions of QCD must also be predictions of the theory in Section 2.3.

The spontaneous symmetry breaking involved in the theory of Section 2.3 should be distinguished from more controversial uses of the term “spontaneous symmetry breaking”. In particular, note that the spontaneously broken $U(N_f) \times U(N_f)$ symmetry is unrelated to the $SU(N_c)$ group in whose Lie algebra the Yang–Mills field is valued. It is sometimes said that the masses of the weak bosons are due to spontaneous symmetry breaking. These masses are due to terms of the form $\phi^\dagger \phi \text{tr}(A_\mu A^\mu)$, with ϕ the Higgs field. When ϕ is expanded about its vacuum expectation value the lowest-order term is an effective mass term for A_μ and the remaining terms soak up the violations of gauge invariance that a bare mass term would incur. By imposing gauge conditions in a particular way you can make this expansion look like spontaneous breaking of gauge transformations. This interpretation poses a problem for the eliminative view of gauge: if the masses of the weak boson were due to the spontaneously broken gauge transformations then removing the gauge structure from the theory would also remove the masses of the weak bosons. However, as Healey (2007, §6.5) argues and Struyve (2011) shows in some detail, you needn’t interpret the Higgs mechanism as a case of spontaneous symmetry breaking. Indeed, like Healey, I think that you shouldn’t. But there’s no need to sort out the correct interpretation of the Higgs mechanism here. The spontaneously broken approximate $U(N_f) \times U(N_f)$ symmetry in the theory of Section 2.3 is not the kind of symmetry for which the Higgs mechanism could be relevant on any interpretation of the latter. The eliminative view of

gauge only motivates one objection to the developments of Section 2.3, and this objection concerns the Yang–Mills vacuum term. But if we treat this term any differently we get bad predictions.

4.1 The chiral anomaly

On the first reading of Healey’s objection, we should remove the Yang–Mills vacuum term from the QCD Lagrangian of Section 2.1. But the presence of such a term isn’t coordinate-invariant: if we try to exclude the vacuum Yang–Mills term from the Lagrangian then it will just come right back when we perform field redefinitions to turn the Higgs–quark coupling terms into mass terms for the quarks. As we saw in Section 2.2, the only coordinate-invariant quantity is $\bar{\theta} = \theta_{\text{YM}} - \theta_{\text{Q}}$, and we can’t assume that θ_{Q} vanishes because there’s CP violation in the weak sector. So the vacuum Yang–Mills term will only stay gone if the reasoning in Section 2.2 fails. But if θ_{YM} and θ_{Q} aren’t interchangeable in this way then the theory predicts the wrong meson decay widths.

The fungibility of θ_{YM} and θ_{Q} amounts to the quantum non-conservation of the Noether current associated with chiral rotations. Any field redefinition is associated with a transformation of the action, and when this transformation is a symmetry of the action it is associated with a current. For a simple analogy, consider a theory of a single scalar field ϕ with action $S = \int d^4x \mathcal{L}$. Under an infinitesimal symmetry transformation the action changes as

$$\phi \mapsto \phi + \epsilon \delta\phi \quad S \mapsto S - \int d^4x \epsilon \partial_\mu j^\mu$$

for some vector field j^μ . This reasoning can be extended to non-compact manifolds by replacing j^μ with a differential form or by imposing boundary conditions, as discussed in Section 3.1.

Any field redefinition corresponds to a transformation of the fields, and when this transformation is a symmetry the field redefinition is associated with a Noether current. Hypotheses about the path integral measure give constraints on the divergence of this current in the quantum theory. In particular, suppose that

$$\mathcal{D}(\phi + \epsilon \delta\phi) = \mathcal{D}\phi \exp\left(i \int d^4x \epsilon \mathcal{A}\right)$$

for some integrand \mathcal{A} . Making the field redefinition that replaces ϕ with $\phi + \epsilon \delta\phi$ then gives

$$Z = \int \mathcal{D}\phi e^{iS} = \int \mathcal{D}\phi e^{iS} \exp\left(-i \int d^4x \epsilon (\partial_\mu j^\mu - \mathcal{A})\right)$$

where j^μ is a Noether current for the symmetry transformation corresponding to this field redefinition. It follows that

$$0 = \frac{i}{Z} \frac{\delta Z}{\delta \epsilon} \Big|_{\epsilon=0} = \frac{\int \mathcal{D}\phi e^{iS} (\partial_\mu j^\mu - \mathcal{A})}{\int \mathcal{D}\phi e^{iS}} = \langle \partial_\mu j^\mu - \mathcal{A} \rangle$$

using the defining property of the generating function Z and the fact that a field redefinition makes no difference to the path integral. This equation is the anomalous Ward–Takahashi identity corresponding to the infinitesimal field

redefinition under consideration, with \mathcal{A} the anomaly. If the anomaly vanishes then the current j^μ is conserved at the quantum level, otherwise the expectation value of its divergence is the expectation value of the anomaly. The reasoning in this paragraph can also be run backward to go from the anomalous Ward–Takahashi identity to the behavior of the path integral measure under field redefinitions.

To show that θ_{YM} and θ_{Q} are interchangeable it thus suffices to show that the anomaly associated with chiral rotations is

$$\mathcal{A}_{\text{chiral}} = -\frac{g^2 N_f}{8\pi^2 N_c} \text{tr}(F_{\mu\nu} \tilde{F}^{\mu\nu})$$

By the reasoning of the previous paragraph, an anomaly of this form is equivalent to a shift in the Yang–Mills vacuum term under a field redefinition. Rephrasing the question this way has two payoffs. First, we can set aside the issue of how to integrate the chiral anomaly over all of Minkowski space. Regardless of what we say about Healey’s second diagnosis, the parameters θ_{YM} and θ_{Q} will be interchangeable if the chiral anomaly has this form. The second payoff is our ability to experimentally determine the expectation value of the chiral rotation current’s divergence.

Various meson decay widths are functions of anomalies, giving the experimental constraint we need. For the simplest example, take the case of neutral pion decay. The primary decay channel for the neutral pion produces two photons. If we assume the picture of Section 2 then we can quickly derive the correct width for this decay. Extend the chiral Lagrangian of Section 2.3 with the lowest-order term compatible with the symmetries and containing the Nambu–Goldstone bosons U_2 , the electromagnetic field $F_{\mu\nu}$, and the volume element $\epsilon^{\mu\nu\alpha\beta}$ (Scherer and Schindler, 2012, §3.5.3). To lowest order in the neutral pion field, this adds a term

$$-\frac{e^2}{16\pi^2 F_\pi} \pi^0 F_{\mu\nu} \tilde{F}^{\mu\nu}$$

The corresponding decay width is

$$\Gamma(\pi^0 \rightarrow \gamma\gamma) = \frac{m_\pi^3}{64\pi} \frac{e^4}{16\pi^4 F_\pi^2} = 7.786 \text{ eV}$$

which is not far from the observed value

$$\Gamma_{\text{obs}}(\pi^0 \rightarrow \gamma\gamma) = 7.63 \text{ eV}$$

This suggests that the majority of the neutral pion’s diphoton decay is due to its coupling with the density $F_{\mu\nu} \tilde{F}^{\mu\nu}$.

Without the chiral anomaly this decay would be suppressed. If the anomaly vanished then chiral rotations would become an approximate symmetry of the effective Lagrangian and thus an approximate symmetry of the effective pion theory. In the effective theory a chiral rotation by ϵ linearly shifts π^0 by ϵ , meaning that a direct coupling between π^0 and $F_{\mu\nu} \tilde{F}^{\mu\nu}$ is disallowed in the massless limit, where chiral rotations are exact symmetries. Such a term could only appear when the masses are turned on and this shift symmetry is explicitly broken, so the lowest order pion–photon coupling must contain m_π^2 . This suppresses the decay width, bringing it down to something on the order of

$$\Gamma_{\mathcal{A}=0}(\pi^0 \rightarrow \gamma\gamma) \lesssim \frac{m_\pi^3}{64\pi} \frac{e^4}{16\pi^4 F_\pi^2} \left(\frac{m_\pi^2}{(4\pi F_\pi)^2} \right)^2 = 1.44 \times 10^{-3} \text{ eV}$$

which badly disagrees with observation. This should bolster our confidence in the ability to trade off between θ_{YM} and θ_{Q} , since this trade-off is enabled by the chiral anomaly and the latter is needed to get the correct neutral pion decay width.

If you reject the size distinction for Healey’s reasons then you cannot reproduce the chiral anomaly, so you will predict the wrong rate for neutral pion decay. There is no way to derive the chiral anomaly in a theory with no gauge structure. Healey does not think that this would be a problem. On his view,

[t]he [chiral] anomaly that permits the two-photon decay of the π^0 -meson involves a violation of gauge symmetry in the “quantum action” $W[A_\mu]$ associated with an external Schwinger source field A_μ . But this source field merely figures as a calculational device for evaluating quantities like vacuum-to-vacuum transition probabilities that are gauge-invariant even though neither their amplitudes nor their generating function $Z[A_\mu] = \exp iW[A_\mu]$ are invariant under “local” gauge transformations in A_μ One can acknowledge the occurrence of anomalies while maintaining that “local” gauge symmetry is a purely formal feature of a theory (Healey, 2007, 183).

The thought seems to be that the derivation of the chiral anomaly can make free use of the gauge structure of the theory without being committed to the physical significance of this structure.

There are a few difficulties with this passage. First, the chiral anomaly involved in neutral pion decay does not involve any violation of gauge invariance. Perhaps the idea is that the chiral anomaly threatens gauge invariance because the integral

$$\int_{\mathbb{M}^4} d^4x \mathcal{A} = -\frac{e^2}{16\pi^2} \int_{\mathbb{M}^4} d^4x F_{\mu\nu} \tilde{F}^{\mu\nu}$$

varies under large gauge transformations. But it *doesn't* vary under large gauge transformations. This was the solution to the “electromagnetic CP problem”: in the case of electromagnetism, where $F_{\mu\nu}$ is a $u(1)$ -valued form, this integral always vanishes. The chiral anomaly does violate a physical symmetry of the classical action upon quantization, since chiral rotations are an approximate symmetry of the classical Lagrangian. But this is an approximate $U(N_f) \times U(N_f)$ phase symmetry that obtains in the limit of vanishing quark mass and is unrelated to whatever Yang–Mills gauge fields are contained in the theory. The approximate phase symmetry is confined to the matter sector, while the gauge group is a feature of both the Yang–Mills sector and the matter sector. The spontaneously broken symmetry in the matter sector is an approximate physical symmetry that rotates quark phases, while gauge transformations are equivalences that mix colors.

Second, the anomaly does not depend on the generating function and sources of the effective action formalism. The anomalous divergence of the chiral current is a claim about the regularized expectation value of an operator, and any way of computing this regularized expectation value is as good as any other. For example, the anomaly is often computed using a “point splitting” method. In the case of neutral pion decay, the anomalous current is given at the classical level by

$$j_3^{5\mu} = \bar{q} \gamma^\mu \gamma^5 \sigma^3 q$$

with q the up–down quark doublet and σ^3 the diagonal Pauli matrix. If we replace the quarks with operator-valued distributions then this product is ill-defined, because products of distributions are only well-defined when the distributions have disjoint singular support. To regularize this product in a gauge-invariant way, consider the “point split” current

$$\epsilon J_3^{5\mu} = \bar{q}(x + \epsilon) \gamma^\mu \gamma^5 \sigma^3 \exp\left(i \int_{x-\epsilon}^{x+\epsilon} A_\mu dx^\mu\right) Q q(x - \epsilon)$$

where A_μ is the electromagnetic gauge potential and Q is a matrix of quark charges. The nonintegrable phase factor is included to maintain gauge invariance: a gauge transformation h sends $q(x)$ to $h(x)q(x)$, so without the phase factor h would lead to an insertion of the nontrivial $h^{-1}(x + \epsilon)h(x - \epsilon)$, breaking gauge invariance. The regulated divergence of this current is

$$\lim_{\epsilon \rightarrow 0} \partial_\mu (\epsilon J_3^{5\mu}) = -\frac{e^2}{16\pi^2} F_{\mu\nu} \tilde{F}^{\mu\nu}$$

reproducing the chiral anomaly (Peskin and Schroeder, 1995, §19.1). The anomaly can be computed without the effective action formalism, and it still does not violate gauge invariance. Indeed, the point-splitting regularization method is Schwinger’s implementation of the idea that “the extraction of gauge-invariant results from a formally gauge-invariant theory is ensured if one employs methods of solution that involve only gauge covariant quantities” (1951, 664).

Finally, Healey’s instrumentalist sentiments in the rest of the passage are difficult to square with his more general project. The question of whether vacuum-to-vacuum transition probabilities are invariant under local gauge transformations is exactly what is at issue. If you accept the size distinction they are not: applying a large gauge transformation changes the path integral measure, thereby changing—for example—the neutron’s EDM. If you reject the size distinction then large gauge transformations preserve the predictions of the theory, and this is the core of Healey’s response to the strong CP problem. But even if the size distinction should be rejected, this isn’t enough to establish Healey’s main claim that gauge transformations have no physical significance. If computations of vacuum-to-vacuum transition amplitudes require a representation that involves gauge transformations then these transformations are significant for physics. Showing that gauge equivalences are a merely mathematical feature of Yang–Mills theory means showing that they can be done away with. You can only acknowledge the chiral anomaly without acknowledging the physical significance of gauge equivalences if you can produce the former without the latter, and this has not been done.

Observations of neutral pion decay show that the current associated with neutral pions has an anomalous divergence. This anomaly is equivalent to the non-invariance of the path integral measure under field redefinitions. It follows that the parameter θ_{YM} appearing in the Yang–Mills vacuum term is coordinate-dependent: by a mere reparametrization of the classical fields we can make the value of θ_{YM} anything we like. So Healey’s view cannot show that θ_{YM} vanishes; this statement has no coordinate-free meaning.

4.2 Masses of mesons

On the second reading of Healey’s objection we can keep the chiral anomaly and thus the correct decay widths. As I noted in Section 2.2, the vacuum Yang–Mills term can make no difference to the action in the case of electromagnetism, because its integral over Minkowski space vanishes. You might think that the same reasoning applies in the case of the strong force. In fact, this *is* what physicists thought in the early days of low-energy hadron physics. If this were right, then the action would always be equal to one in which θ_{YM} vanishes, no matter the value of θ_{YM} in the Lagrangian. But it conflicts with experiment. If the integral vanishes then there must be a pseudoscalar meson with the quantum numbers of the η' and the mass of the pions, and there isn’t one.

If the vacuum Yang–Mills term of the action vanishes for all values of θ_{YM} then the symmetry-breaking pattern of the QCD action is not the one described in Section 2.3. The QCD Lagrangian has an approximate $U(2) \times U(2)$ symmetry group, since the up and down quarks are very light. Whether this is also the symmetry group of the quantum effective action depends on our attitude about the Yang–Mills vacuum term. In particular, a chiral rotation

$$q \mapsto q + i\epsilon\gamma^5 q$$

is accompanied by a phase rotation on the path integral measure

$$\mathcal{D}A \mathcal{D}q \mathcal{D}\bar{q} \mapsto \mathcal{D}A \mathcal{D}q \mathcal{D}\bar{q} \exp\left(-i\frac{\epsilon g^2}{8\pi^2 N_c} \int d^4x \operatorname{tr}(F_{\mu\nu} \tilde{F}^{\mu\nu})\right)$$

In Section 2.3 we assumed that the integral in this phase does not vanish, and that therefore the path integral measure changes under chiral rotations. As such, we found only three Nambu–Goldstone bosons, corresponding to the remaining three broken generators of $U(2) \times U(2)$. However, if we suppose that the integral does vanish then the every element of $U(2) \times U(2)$ gives a symmetry of the effective action. The chiral condensate then leads to the spontaneous breaking of four approximate symmetry generators and four accompanying Nambu–Goldstone bosons.

This alternative pattern of spontaneous symmetry breaking means we have to revise the effective chiral Lagrangian of Section 2.3. By rotating $\bar{\theta}$ into the Yang–Mills vacuum term we can eliminate it, giving the chiral Lagrangian

$$\mathcal{L}_{\chi,2} = \frac{1}{4}F_\pi^2 \operatorname{tr}(\partial^\mu \tilde{U}_2^\dagger \partial_\mu \tilde{U}_2) + \frac{1}{2}B_\pi F_\pi^2 \operatorname{tr}(M\tilde{U}_2 + M\tilde{U}_2^\dagger)$$

where \tilde{U}_2 is $U(2)$ -valued, rather than $SU(2)$ -valued, taking the form

$$\tilde{U}_2 = \exp\left(\frac{i}{F_\pi} \begin{pmatrix} \eta_0 + \eta_3 & \sqrt{2}\pi^- \\ \sqrt{2}\pi^+ & \eta_0 - \eta_3 \end{pmatrix}\right)$$

Since chiral rotations are now a symmetry of theory we can eliminate the phase on the mass matrix by absorbing it into \tilde{U}_2 . As before, we can determine $F_\pi = 92.07 \text{ MeV}$ from the leptonic decay rate of the pion and interpret B_π in terms of the vacuum energy. The new field η_0 corresponds to the broken generator of chiral rotations. Because this generator is proportional to the identity matrix on flavor space we expect the neutral pion and the newly predicted particle to be some mixture of η_0 and η_3 that depends on the up–down mass difference.

This Lagrangian is not a good model of low-energy QCD. It fails at even the most basic task of getting the meson masses right. Expanding about the minimum of \tilde{U}_2 and assuming that the up–down flavor symmetry is exact, we find the three pions along with a new particle η , all of which have the mass

$$m_{\pi^0}^2 = m_{\pi^\pm}^2 = m_\eta^2 = B_\pi(m_u + m_d)$$

So the new Nambu–Goldstone boson must have a mass comparable to the pions’. The η has the right quantum numbers to be the new particle, but it has a mass of 575 MeV—far heavier than the 140 MeV pions. Including the mass difference between the up and down quarks does not help; this will cause some mixing between η_0 and η_3 , but no amount of mixing can raise the η ’s predicted mass to 575 MeV. So there are two possibilities: either this theory gets the mass of the η badly wrong, or it predicts a heretofore undiscovered pseudoscalar meson whose mass is comparable to the pions’. Since a particle like this could not have escaped our notice, the theory has a serious problem either way.

An obvious fix suggests itself: the problem might be avoided if we incorporate contributions from heavier quarks. Heavier quarks certainly contribute something to the dynamics of the pions. According to the quark model, a pseudoscalar meson like the charged pion π^+ is a bound state of a quark and an antiquark—in the case of π^+ , an up quark and an antidown quark. One expression of this fact is that a meson’s quantum numbers are determined by those of its constituent quarks. For example, the π^+ has unit charge, the sum of the up quark’s $2/3$ charge and antidown’s $1/3$ charge. But really the up and the antidown are just the valence quarks of the π^+ . The charged pion is an effective low-energy degree of freedom that represents fluctuations in the chiral condensate; it has contributions from every quark. The heavier quarks just contribute less than the up and down do, because the pion lives at low energies. Because the η is observed to be heavier, we might expect much of its mass to be due to contributions from heavier quarks.

We can include the effects of the strange quark by assuming that it’s approximately massless, too. On this assumption, the reasoning in Section 2.3 goes through as before, except this time the QCD Lagrangian has a $U(3) \times U(3)$ chiral flavor-mixing symmetry that’s spontaneously broken to its diagonal $U(3)$. It’s much less plausible to treat the strange quark as massless than it is to treat the up and down as massless—while the up and down quarks have masses of only a few MeV, the strange quark is closer to 100 MeV. But the proof of the pudding is in the eating, and it turns out that the approximation works in practice.

We do it the right way first (Scherer and Schindler, 2012, Ch. 3). Since the chiral anomaly explicitly breaks the chiral rotation symmetry at the quantum level there are eight Nambu–Goldstone bosons. We take these to be the pions, kaons, and the η . So consider the $SU(3)$ -valued field

$$U_3 = \exp \left(\frac{i}{F_0} \begin{pmatrix} \eta_3 + \frac{1}{\sqrt{3}}\eta_8 & \sqrt{2}\pi^+ & \sqrt{2}K^+ \\ \sqrt{2}\pi^- & -\eta_3 + \frac{1}{\sqrt{3}}\eta_8 & \sqrt{2}K^0 \\ \sqrt{2}K^- & \sqrt{2}K^0 & -\frac{2}{\sqrt{3}}\eta_8 \end{pmatrix} \right)$$

whose effective Lagrangian is, to lowest order,

$$\mathcal{L}_{\chi,3} = \frac{1}{4}F_0^2 \text{tr} \left(\partial^\mu U_3^\dagger \partial_\mu U_3 \right) + \frac{1}{2}B_0F_0^2 \text{tr} \left(e^{-i\bar{\theta}/3} M U_3 + e^{i\bar{\theta}/3} M U_3^\dagger \right)$$

We can determine the decay constant F_0 and the vacuum parameter B_0 as we did in the two-flavor case. The subscript on η_8 signifies that it is the Nambu–Goldstone boson associated with the eighth Gell-Mann matrix. Mass differences between the quarks make the neutral pion and the η meson mixtures of η_8 and η_3 .

This theory successfully accounts for the masses of the eight lightest pseudoscalar mesons in Table 1. Expanding about the minimum of U_3 gives, to zeroth order in the up–down mass difference, the mass terms

$$\begin{aligned} m_\pi^2 &= B_0(m_u + m_d) \\ m_{K^\pm}^2 &= B_0(m_u + m_s) & m_{K^0}^2 &= B_0(m_d + m_s) \\ m_\eta^2 &= \frac{1}{3}B_0(m_u + m_d + 4m_s) \end{aligned}$$

There are three light pions, then four kaons of roughly the same mass, then one η whose mass is slightly larger than the kaons'. Once F_0 and B_0 are fixed, this model gives good predictions for various semileptonic kaon decays, pion–kaon scattering, and electromagnetic form factors, in addition to reproducing the predictions of the $SU(2) \times SU(2)$ theory.⁹

Now suppose that the Yang–Mills vacuum term vanishes; the resulting theory again predicts a nonexistent light particle (Weinberg, 1975). The path integral measure becomes invariant under chiral rotations, Goldstone’s theorem applies, and a new field η_0 appears. The Nambu–Goldstone fields now assemble into a $U(3)$ matrix, making the field in the effective theory

$$\tilde{U}_3 = \exp\left(\frac{i\sqrt{2}}{\sqrt{3}F_0}\eta_0\right)U_3$$

Rotating $\bar{\theta}$ into the vanishing Yang–Mills vacuum term and expanding about the minimum of \tilde{U}_3 gives the same masses as before for the charged pions and the kaons, but the other three fields are no longer mass eigenstates. The η_3 – η_8 – η_0 mass matrix is

$$M_\eta^2 = B_0 \begin{pmatrix} m_u + m_d & \frac{1}{\sqrt{3}}(m_u - m_d) & \frac{\sqrt{2}}{\sqrt{3}}(m_u - m_d) \\ \frac{1}{\sqrt{3}}(m_u - m_d) & \frac{1}{3}(m_u + m_d + 4m_s) & \frac{\sqrt{2}}{3}(m_u + m_d - 2m_s) \\ \frac{\sqrt{2}}{\sqrt{3}}(m_u - m_d) & \frac{\sqrt{2}}{3}(m_u + m_d - 2m_s) & \frac{2}{3}(m_u + m_d + m_s) \end{pmatrix}$$

in the (η_3, η_8, η_0) basis. The strange quark is much heavier than the other two, so to a first approximation we can assume that m_u and m_d vanish, giving two massless eigenstates and one with mass

$$m_\eta^2 = 2B_0m_s$$

Projecting the mass matrix onto the eigenspace spanned by the other two modes gives

$$M_{\eta_3, \eta_0}^2 = B_0 \begin{pmatrix} m_u + m_d & \frac{\sqrt{3}}{\sqrt{2}}(m_u - m_d) \\ \frac{\sqrt{3}}{\sqrt{2}}(m_u - m_d) & \frac{3}{2}(m_u + m_d) \end{pmatrix}$$

⁹See Donoghue et al. (1992, Ch. VII) and Scherer and Schindler (2012, §5.3.6) for further discussion of these results.

so in the limit where the up and down have identical non-vanishing masses we have two more mass eigenstates

$$m_{\pi_0}^2 = B_0(m_u + m_d) \quad m_{\eta'}^2 = \frac{3}{2}B_0(m_u + m_d)$$

Whether or not the Yang–Mills vacuum term vanishes we obtain the same masses for the pions and kaons, and the mass of the η will be roughly the same in both cases, as well. But when the vacuum Yang–Mills term vanishes there must be a ninth pseudoscalar meson with mass $\sqrt{3/2}m_\pi$, and there is no such meson. More careful treatments of the η_3 – η_8 – η_0 mixing can raise the value of $m_{\eta'}$ as high as $\sqrt{3}m_\pi$, but this still isn't enough. When the Yang–Mills vacuum term vanishes there just isn't enough mass to go around.

The missing η' mass lives in the Yang–Mills vacuum term (Veneziano, 1979; Witten, 1979, 1980). In the limit of vanishing quark mass chiral rotations do not become symmetries of the effective action, since they change the path integral measure. However, they do become symmetries in the limit where the quark masses and N_c^{-1} both vanish—that is, when the number of colors is large—because the chiral anomaly is proportional to N_c^{-1} . We therefore suppose that η_0 receives a mass contribution of order N_c^{-1} , giving the effective Lagrangian

$$\begin{aligned} \mathcal{L}_{\chi, \bar{3}} = & \frac{1}{4}F_0^2 \text{tr}(\partial^\mu \tilde{U}_3^\dagger \partial_\mu \tilde{U}_3) + \frac{1}{2}B_0F_0^2 \text{tr}(e^{-i\bar{\theta}/3}M\tilde{U}_3 + e^{i\bar{\theta}/3}M\tilde{U}_3^\dagger) \\ & + \frac{\chi}{2N_c^2}(\log \det \tilde{U}_3)^2 \end{aligned}$$

Consider the limit in which the up and down masses are the same, so that the masses of the pions and kaons are unchanged and the squared mass matrix is diagonal except for the terms mixing η_8 and η_0 . These terms can be written

$$M_{\eta_8, \eta_0}^2 = \frac{1}{3} \begin{pmatrix} 4m_K^2 - m_\pi^2 & 2\sqrt{2}(m_\pi^2 - m_K^2) \\ 2\sqrt{2}(m_\pi^2 - m_K^2) & 2m_K^2 + m_\pi^2 + \frac{18}{F_0^2 N_c^2} \chi \end{pmatrix}$$

Since the trace of a matrix is the sum of its eigenvalues this implies the Witten–Veneziano relation

$$m_\eta^2 + m_{\eta'}^2 = 2m_K^2 + \frac{6}{F_0^2 N_c^2} \chi$$

In the limit of vanishing quark masses the η and the kaons become massless as well, but the η' still receives nonzero contributions from the chiral anomaly. This explains why the η' is so much heavier than the other mesons in Table 1.

The coefficient χ can be determined by matching the vacuum energy of this effective theory with the vacuum energy of the QCD Lagrangian, just as we did with the coefficient of the mass term (Witten, 1979). If we absorb $\bar{\theta}$ into \tilde{U}_3 by redefining η_0 then we can write this Lagrangian as

$$\begin{aligned} \mathcal{L}_{\chi, \bar{3}} = & \frac{1}{4}F_0^2 \text{tr}(\partial^\mu \tilde{U}_3^\dagger \partial_\mu \tilde{U}_3) + \frac{1}{2}B_0F_0^2 \text{tr}(M\tilde{U}_3 + M\tilde{U}_3^\dagger) \\ & - \frac{\chi}{2N_c^2}(\bar{\theta} - \log \det \tilde{U}_3)^2 \end{aligned}$$

In the large N_c limit the last term is suppressed, so the vacuum expectation value of \tilde{U}_3 is still approximately the identity matrix and the mass term still

accounts for vacuum energy contributions from the quark sector. So in the vacuum the last term is of order $\bar{\theta}^2$ and due entirely to the Yang–Mills sector. Putting $\bar{\theta}$ into the Yang–Mills term and expanding in terms of $\bar{\theta}$ gives

$$\frac{\chi}{N_c^2} = \left(\frac{g^2}{8\pi^2 N_c} \right)^2 \int d^4x \left\langle \text{tr} \left(F_{\mu\nu}(x) \tilde{F}^{\mu\nu}(x) \right) \text{tr} \left(F_{\alpha\beta}(0) \tilde{F}^{\alpha\beta}(0) \right) \right\rangle$$

Lattice computations of this expectation value are consistent with the Witten–Veneziano relation (ETM Collaboration, 2015).

A theory in which the Yang–Mills vacuum term vanishes predicts a nonexistent particle, or perhaps the wrong mass for the η or η' . It's hard to see how the eliminative view can avoid this fate. Healey doesn't address the issue. However, his discussion of the strong CP problem in path-based formulations of Yang–Mills theory refers to Fort and Gambini (2000) as demonstrating that the strong CP problem is solved in a loop-based formulation, so we might look to them for a response. But we won't find one. Healey aims to interpret Yang–Mills theory without referring to gauge transformations at all. This is what makes his solution to the strong CP problem tick: if we can't distinguish large and small gauge transformations then we can't introduce the Yang–Mills vacuum term, or its integral must always vanish, and so there can be no observable CP violation in the strong sector. But Fort and Gambini's formalism includes the same gauge structure as the usual field formalism. Their results are also unresponsive to the problem of the η' mass.

Fort and Gambini introduce gauge transformations to allow for interactions between the Yang–Mills configuration and the matter fields. As they say, in an analysis of the η' mass

the fermionic degrees of freedom must be included in the holonomy formulation. This was done some years ago by including gauge-invariant hadronic objects built on open paths, in addition to the closed ones or loops for the pure gauge theory, giving rise to the so-called *P-representation* The P is for paths, which, in this case when matter fields are present, are in general open. In what follows, although strictly holonomies are defined for closed paths, we will use the term holonomy representation indistinguishably from P-representation (2000, 344).

In Fort and Gambini's P-representation the Yang–Mills configuration is represented by a map that assigns an element of $SU(N_c)$ to every path in spacetime. An assignment H couples to the quark fields through terms of the form

$$\bar{q}(y) H(\gamma) q(x)$$

where γ is a curve from x to y . Like gauge potentials, such maps contain “too much” information: different maps H and H' assigning group elements to curves in spacetime correspond to the same physical state of affairs if there is a smooth function $h : \mathbb{M}^4 \rightarrow SU(N_c)$ such that

$$H'(\gamma) = h(y) \cdot H(\gamma) \cdot h^{-1}(x)$$

for all points x and y in \mathbb{M}^4 and all curves γ from x to y . In other words, two maps H and H' are equivalent just in case they are related by a gauge

transformation. Since there is a bijection between $\mathfrak{su}(N_c)$ -valued one-forms and smooth maps assigning elements of $SU(N_c)$ to open curves in spacetime, the P-representation is just a notational variant of the gauge potential formalism (Schreiber and Waldorf, 2009).

Since the P-representation has the same gauge structure as the gauge potential formalism, it isn't of any use to Healey. Nor do Fort and Gambini's more informal remarks offer any help. They claim first that the problem of the η' mass is solved by the gauge invariance of the theory:

the anomaly occurs as a consequence of the incompatibility of two classical symmetries—gauge and chiral invariance—at the quantum level. It happens that the gauge symmetry may only be preserved at the price of sacrificing the chiral symmetry which become anomalously broken. The P-representation deals with gauge-invariant quantities and hence has no chance to implement the chiral symmetry (2000, 345).

It's true that gauge transformations are incompatible with chiral rotations in the quantum theory. This can be seen in the fact that the path integral measure picks up a phase under a field redefinition, and it can also be articulated using algebraic tools (Strocchi, 2013, §8.2). But the P-representation and the gauge potential formalism are no different with respect to this incompatibility. The two formalisms have the same gauge transformation structure in the Yang–Mills sector, and they are identical in the matter sector. Since chiral rotations are transformations in the matter sector, they can be implemented in the P-representation exactly as above:

$$q \mapsto q + i\epsilon\gamma^5 q$$

At the classical level the P-representation has the same symmetry groups as the gauge potential formalism. Fort and Gambini do not explain why this symmetry structure is broken at the quantum level. They say that the puzzle of the η' mass is solved because the P-representation “does not bear [chiral rotation] symmetry at the second quantized level” (2000, 347), but of course neither does the gauge potential formalism. The only theory that has chiral symmetry at the quantum level is one in which the integral of the vacuum Yang–Mills term does not contribute to the quantum action. And a theory like this gets the facts wrong.

5 Conclusion

Sometimes gauge transformations are physical symmetries; the gauge potentials they relate represent different physical states of affairs. In particular, “large” gauge transformations can change the physical facts. This is a problem for the standard philosophical interpretation of gauge, on which gauge transformations can be eliminated from any gauge theory without loss of physical content. It also engenders a problem for the Standard Model of particle physics. For if large gauge transformations are symmetries and not redundancies then the strong force violates CP—unless a parameter from the Yang–Mills sector and a parameter from the quark sector conspire to prevent it. Healey argues that

we should avoid the latter problem by denying that large gauge transformations are symmetries. I have argued above that this isn't a viable strategy: there is no way to implement this suggestion without running afoul of well-confirmed features of particle physics.

I have argued that the eliminative view of gauge gives incorrect predictions about mesons. On this view gauge transformations are a convenience of one particular formulation of Yang–Mills theory, and we would have a truer representation of the facts if we were to do away with them entirely. Section 3 showed that if the eliminative view were true then the vacuum Yang–Mills term

$$\frac{g^2\theta_{\text{YM}}}{16\pi^2 N_c} \int_K d^4x \operatorname{tr}(F_{\mu\nu}\tilde{F}^{\mu\nu})$$

would lead to inconsistency when integrated over any region K . By Stokes' theorem it is a matter of mathematical fact that this integral coincides with the integral

$$\frac{g^2\theta_{\text{YM}}}{16\pi^2 N_c} \int_{\partial K} d^3x \epsilon^{\mu\nu\alpha} \operatorname{tr}\left(A_\mu F_{\nu\alpha} - i\frac{2g}{3\sqrt{N_c}} A_\mu A_\nu A_\alpha\right)$$

But this integral varies under large gauge transformations. So if we were to eliminate gauge from the theory then each configuration would be assigned contradictory values for the vacuum Yang–Mills term of the action: one for each class of representative gauge potentials that differ by a large gauge transformation. As I argued in Section 4, any attempt to excise this inconsistency from the Standard Model leads to a bad prediction of the decay widths or masses of some mesons: making this term vanish by setting $\theta_{\text{YM}} = 0$ gives the wrong decay widths, and taking the integral to vanish gives the wrong masses. There is no obvious way for the eliminative view to reproduce the effects of the vacuum Yang–Mills term, and I argued that attempted solutions following Healey's suggestions will not work. Healey is right to say that the eliminative view dissolves the strong CP problem outlined in Section 2, but this dissolution comes at too high a cost.

The results above tell us that we should reject the eliminative view, but they don't tell us much more. There are good reasons behind the eliminative view. For one thing, quantization procedures for and computations in gauge field theories treat gauge-equivalent configurations as the same physical state. For another, gauge-equivalent configurations are empirically and dynamically indistinguishable—you and I aren't able to observe any differences between two gauge-equivalent potentials and the dynamics can't choose between two gauge-equivalent possible futures for some initial state. So it would be too hasty to conclude that all gauge transformations are symmetries. But it's also not enough to simply make an exception for large gauge transformations. Do we make an exception for any gauge transformation that's nontrivial on the boundary of any region? Only those on the sphere at infinity that also spoil the gauge invariance of the vacuum Yang–Mills term? Something in between? What other consequences does this exception have, if any? Why think it's consistent with the rest of the theory? Why think that this exception will be sufficient—in the case of QCD it saves some meson phenomena, but what about other gauge theories and other phenomena? The answers to these questions should follow from a general understanding of gauge structure that implies an appropriate distinction between large and small gauge transformations.

Let me end by suggesting one such general understanding. The distinction between large and small gauge transformations is an instance of the fact, recently stressed by Belot (2018), that physicists will take isomorphic models of a theory to be distinct physical states of affairs in some theories and in some contexts. They will sometimes do so even when the transformation is spacetime-dependent, as above and as in monopole physics and gravitation. Like Belot, I think we ought to understand this attitude by looking to structures on the set of models—in the case of Yang–Mills theory, the set of gauge configurations. In particular, we should understand the set of gauge transformations between two configurations to be part of the structure postulated by the theory. Gauge transformations tell us when two configurations represent the same physical state of affairs, but they tell us more than this. In electromagnetism the set of gauge transformations over \mathbb{M}^4 from A_μ to A'_μ is the set of smooth functions $h : \mathbb{M}^4 \rightarrow U(1)$ such that

$$A'_\mu = A_\mu - \frac{i}{g} h \partial_\mu h^{-1}$$

It would be a different theory if we instead took the gauge transformations to be smooth functions $\lambda : \mathbb{M}^4 \rightarrow \mathbb{R}$ such that

$$A'_\mu = A_\mu + \frac{1}{g} \partial_\mu \lambda$$

It would be different because it counts gauge transformations differently. The eliminative view of gauge collapses these two theories into one, thereby losing part of the theory’s structure. But it’s this structure that governs, for example, the gauge-invariant imposition of boundary conditions in Section 3.1. Nontrivial gauge structure naturally gives rise to a distinction between large and small gauge transformations, and eliminating the gauge structure also eliminates the distinction between large and small. But these details have to wait for another day.

To bring the Lagrangian into a form with no CP violation in the matter sector we must perform the transformation

$$q \mapsto e^{i\theta_Q \gamma^5 / 2N_f} \mathbf{1} q$$

where we have inserted an explicit $N_f \times N_f$ identity matrix $\mathbf{1}$. Assuming that $(-\theta_Q)$ is infinitesimal, we have

$$\det \mathcal{U} = \exp \operatorname{tr} \log \mathcal{U} = \exp \left(i \frac{\theta_Q}{2} \int d^4 x \operatorname{tr} (\gamma^5) \delta^4(x-x) \right)$$

so

$$\begin{aligned} \mathcal{D}q \mathcal{D}\bar{q} &\mapsto \mathcal{D}q \mathcal{D}\bar{q} \exp \left(-i\theta_Q \int d^4 x \operatorname{tr} (\gamma^5) \delta^4(x-x) \right) \\ &= \mathcal{D}q \mathcal{D}\bar{q} \exp \left(i \frac{1}{2N_f} \int d^4 x \theta_Q \mathcal{A}(x) \right) \end{aligned}$$

Now to compute $\mathcal{A}(x)$. We write

$$\begin{aligned}
\mathcal{A}(x) &= \lim_{y \rightarrow x} -2 \operatorname{tr}(\gamma^5 \mathbf{1} f(-\mathcal{D}_x^2/M^2)) \delta^4(x-y) \\
&= -2 \lim_{y \rightarrow x} \int \frac{d^4 k}{(2\pi)^4} \operatorname{tr}(\gamma^5 \mathbf{1} f(-\mathcal{D}_x^2/M^2)) e^{ik \cdot (x-y)} \\
&= -2M^4 \lim_{y \rightarrow x} \int \frac{d^4 k}{(2\pi)^4} \operatorname{tr}(\gamma^5 \mathbf{1} f(-[ik + \mathcal{D}_x/M]^2)) \\
&= - \int \frac{d^4 k}{(2\pi)^4} f''(k^2) \operatorname{tr}(\gamma^5 \mathbf{1} \mathcal{D}_x^4)
\end{aligned}$$

Now

$$\begin{aligned}
\mathcal{D}_x^2 &= \frac{1}{4} \{(D_x)^\mu, (D_x)^\nu\} \{\gamma_\mu, \gamma_\nu\} + \frac{1}{4} [(D_x)^\mu, (D_x)^\nu] [\gamma_\mu, \gamma_\nu] \\
&= \frac{1}{4} \{(D_x)^\mu, (D_x)^\nu\} \{\gamma_\mu, \gamma_\nu\} + \frac{1}{4} [(D_x)^\mu, (D_x)^\nu] [\gamma_\mu, \gamma_\nu] \\
&= D_x^2 + i \frac{g}{4\sqrt{N_c}} F_{\mu\nu} [\gamma_\mu, \gamma_\nu]
\end{aligned}$$

So

$$\begin{aligned}
\mathcal{A}(x) &= \frac{g^2 N_f}{256 N_c \pi^4} \int d^4 k f''(k^2) \operatorname{tr}(F_{\mu\nu} F_{\alpha\beta}) \operatorname{tr}(\gamma^5 [\gamma^\mu, \gamma^\nu] [\gamma^\alpha, \gamma^\beta]) \\
&= - \frac{g^2 N_f}{8 N_c \pi^2} \operatorname{tr}(F_{\mu\nu} \tilde{F}^{\mu\nu})
\end{aligned}$$

Meaning, finally, that

$$\mathcal{D}q \mathcal{D}\bar{q} \mapsto \mathcal{D}q \mathcal{D}\bar{q} \exp\left(-i \frac{\theta_Q g^2}{16\pi^2 N_c} \int d^4 x \operatorname{tr}(F_{\mu\nu} \tilde{F}^{\mu\nu})\right)$$

References

- Adler, S. L. (1969). Axial-vector vertex in spinor electrodynamics. *Physical Review*, 177:2426–2438.
- Anderson, I. M. (1992). Introduction to the variational bicomplex. In Gotay, M., Marsden, J. E., and Moncrief, V. E., editors, *Mathematical Aspects of Classical Field Theory*, pages 51–73. American Mathematical Society.
- Bain, J. (2019). Why be natural? *Foundations of Physics*.
- Barr, S. M. (1984). Solving the strong cp problem without the peccei-quinn symmetry. *Physical Review Letters*, 53(4):329.
- Bell, J. and Jackiw, R. (1969). A PCAC puzzle: $\pi^0 \rightarrow \gamma\gamma$ in the σ -model. *II Nuovo Cimento A*, 60(1):47–61.
- Belot, G. (2018). Fifty million Elvis fans can't be wrong. *Noûs*, 52:946–981.
- Bleecker, D. (1981). *Gauge Theory and Variational Principles*. Addison-Wesley.

- Callan Jr., C. G., Coleman, S., Wess, J., and Zumino, B. (1969). Structure of phenomenological Lagrangians. II. *Physical Review*, 177(5):2247.
- Cichy, K., Garcia-Ramos, E., and Jansen, K. (2013). Chiral condensate from the twisted mass dirac operator spectrum. *Journal of High Energy Physics*, 2013(10):175.
- Coleman, S., Wess, J., and Zumino, B. (1969). Structure of phenomenological Lagrangians. I. *Physical Review*, 177(5):2239.
- Crewther, R., Di Vecchia, P., Veneziano, G., and Witten, E. (1979). Chiral estimate of the electric dipole moment of the neutron in quantum chromodynamics. *Physics Letters B*, 88(1-2):123–127.
- Donoghue, J. F., Golowich, E., and Holstein, B. R. (1992). *Dynamics of the Standard Model*. Cambridge University Press.
- Earman, J. (2004). Curie’s principle and spontaneous symmetry breaking. *International Studies in the Philosophy of Science*, 18(2-3):173–198.
- Ellis, J. and Gaillard, M. K. (1979). Strong and weak CP violation. *Nuclear Physics B*, 150:141–162.
- Fort, H. and Gambini, R. (2000). U(1) puzzle and the strong CP problem from a holonomy perspective. *International Journal of Theoretical Physics*, 39:341–349.
- Fujikawa, K. (1979). Path-integral measure for gauge-invariant fermion theories. *Physical Review Letters*, 42(18):1195.
- Fujikawa, K. (1980). Path integral for gauge theories with fermions. *Physical Review D*, 21(10):2848.
- Healey, R. (2007). *Gauging What’s Real*. Oxford University Press.
- Healey, R. (2010). Gauge symmetry and the theta-vacuum. In Suárez, M., Dorato, M., and Rédei, M., editors, *EPSA Philosophical Issues in the Sciences: Launch of the European Philosophy of Science Association*, pages 105–116. Springer.
- Irastorza, I. G. and Redondo, J. (2018). New experimental approaches in the search for axion-like particles. *Progress in Particle and Nuclear Physics*, 102:89–159.
- Marsden, J. E. and Hughes, T. J. R. (1994). *Mathematical Foundations of Elasticity*. Dover.
- Nelson, A. (1984). Naturally weak CP violation. *Physics Letters B*, 136(5–6):387–391.
- Peccei, R. D. and Quinn, H. R. (1977a). Constraints imposed by CP conservation in the presence of pseudoparticles. *Physical Review D*, 16:1791–1797.
- Peccei, R. D. and Quinn, H. R. (1977b). CP conservation in the presence of pseudoparticles. *Physical Review Letters*, 38:1440–1443.

- Peskin, M. E. and Schroeder, D. V. (1995). *An Introduction to Quantum Field Theory*. Addison-Wesley.
- Rubakov, V. (2002). *Classical Theory of Gauge Fields*. Princeton University Press.
- Saunders, D. J. (1989). *The Geometry of Jet Bundles*. Cambridge University Press.
- Scherer, S. and Schindler, M. R. (2012). *A Primer for Chiral Perturbation Theory*. Springer.
- Schreiber, U. and Waldorf, K. (2009). Parallel transport and functors. *Journal of Homotopy and Related Structures*, 4:187–244.
- Schwinger, J. (1951). On gauge invariance and vacuum polarization. *Physical Review*, 82:664–679.
- Srednicki, M. (2007). *Quantum Field Theory*. Cambridge University Press.
- Strocchi, F. (2013). *An Introduction to Non-Perturbative Foundations of Quantum Field Theory*. Oxford University Press.
- Struyve, W. (2011). Gauge invariant accounts of the Higgs mechanism. *Studies in History and Philosophy of Modern Physics*, 42(4):226–236.
- 't Hooft, G. (1976a). Computation of the quantum effects due to a four-dimensional pseudoparticle. *Physical Review D*, 14:3432–3450.
- 't Hooft, G. (1976b). Symmetry breaking through Bell–Jackiw anomalies. *Physical Review Letters*, 37:8–11.
- 't Hooft, G. (1979). Naturalness, chiral symmetry, and spontaneous chiral symmetry breaking. *NATO Advanced Science Institutes Series B: Physics*, 59:135–157.
- 't Hooft, G. (1986). How instantons solve the $U(1)$ problem. *Physics Reports*, 142(6):357–387.
- 't Hooft, G. and Veltman, M. (1972). Regularization and renormalization of gauge fields. *Nuclear Physics B*, 44(1):189–213.
- Tanabashi, M. et al. (Particle Data Group) (2018). Review of particle physics. *Physical Review D*, 98(3):030001.
- The ETM collaboration, Cichy, K. Garcia-Ramos E. Jansen K. Ottnad K. and Urbach C. (2015). Non-perturbative test of the Witten-Veneziano formula from lattice QCD. *Journal of High Energy Physics*, 2015(9):20.
- Veneziano, G. (1979). $U(1)$ without instantons. *Nuclear Physics B*, 159(1–2):213–224.
- Vicari, E. and Panagopoulos, H. (2009). θ dependence of $SU(N)$ gauge theories in the presence of a topological term. *Physics Reports*, 470(3–4):93–150.
- Weinberg, S. (1975). The $U(1)$ problem. *Physical Review D*, 11(12):3583–3593.

- Weinberg, S. (1995). *The Quantum Theory of Fields*. Cambridge University Press.
- Williams, P. (2015). Naturalness, the autonomy of scales, and the 125 GeV Higgs. *Studies in History and Philosophy of Modern Physics*, 51:82–96.
- Witten, E. (1979). Current algebra theorems for the $U(1)$ “Goldstone boson”. *Nuclear Physics B*, 156:269–283.
- Witten, E. (1980). Large N chiral dynamics. *Annals of Physics*, 128(2):363–375.
- Wu, T. T. and Yang, C. N. (1975). Concept of nonintegrable phase factors and global formulation of gauge fields. *Physical Review D*, 12:3845–3857.