

Contents

0	Overview of the lectures	3
1	Recap: path integrals and Feynman diagrams	4
1.1	Gaussian integration and Wick's theorem	4
1.2	Non-convergence of the perturbative expansion	4
1.3	0+1d free scalar field	5
1.4	ϕ^4 interaction in 0+1d	6
1.5	Aside: Brillouin-Wigner perturbation theory in quantum mechanics	8
1.6	Vacuum energy of a 3+1d free scalar field	11
1.7	The Casimir effect	11
2	The Lagrangian of the Standard Model	14
2.1	Gauge fields	14
2.2	Why gauge invariance?	15
2.3	Coupling gauge fields to scalars	16
2.4	Two-component spinors	18
2.5	Coupling gauge fields to fermions	21
2.6	Coupling scalars to fermions	22
2.7	The Lagrangian of the Standard Model	22
2.8	Questions	25
2.9	Anomaly cancellation	25
3	Anomalies	27
3.1	Aharonov-Bohm phase and the Dirac quantization of monopole charges	27
3.2	2d fermions	29
3.3	An index theorem	32
3.4	Anomalies of 4d chiral fermion	34
4	Spontaneous symmetry breaking, Higgs effect and solitons	35
4.1	Spontaneous symmetry breaking	36
4.2	Higgs mechanism in U(1) gauge theory	37
4.3	Under the external magnetic field	39
4.4	Vortex solution	40
4.5	Bogomolny trick	41
4.6	Higgs mechanism in non-Abelian gauge theory	43
4.7	Monopole solution	46
4.8	Higgs mechanism in the Standard model	49
5	Renormalization group	53
5.1	Scalar ϕ^4	53
5.2	Running of the gauge couplings	56

5.3	Two-loop running and the fixed points	59
5.4	The Banks-Zaks IR fixed point	62
5.5	The Litim-Sannino UV fixed point	63
6	Qualitative discussions of strongly-coupled gauge theories	64
6.1	$N_f = 0$: color confinement and the mass gap	65
6.2	$N_f = 1, 2, \dots$: chiral symmetry breaking and the U(1) problem	66
6.3	The U(1) problem	68
6.4	Chiral Lagrangian	72
6.5	Baryon as a soliton in the chiral Lagrangian	74
6.6	Wess-Zumino-Witten term	75
7	Renormalizability, effective field theory, and UV completeness	78
7.1	Assigning dimensions to operators	78
7.2	Renormalizability in the traditional sense	79
7.3	Non-renormalizable theories as effective theories	80
7.4	‘Completion’ of an effective theory	82
7.5	Standard Model as an effective field theory	84
7.6	Renormalizability of Yang-Mills and gravity in various dimensions	86
	References	88

0 Overview of the lectures

The name of this slot is *Theory of Elementary Particles* but it is used essentially as *Quantum Field Theory III* following *Quantum Field Theory I and II*. The aim of the lectures is to learn various theoretical aspects of quantum field theory (QFT) in the context of the Standard Model. The experimental aspects will be covered by *Elementary Particle Physics, I, II and III*.

One of the difficulties involved in studying QFT is that it not only involves many new conceptual points but also many technically complicated computations. This makes it rather hard for a first learner to distinguish which is which. I would aim to use simpler toy models which emphasize the conceptual points, with the caveat that the conclusions obtained from them cannot usually be compared with experiments. My normalization of various fields does not follow the accepted conventions in the Standard Model. So if you actually do computations for it, please do not use mine.

I will use the unit where $c = \hbar = 1$. The metric is mostly plus, so $x_\mu x^\mu = -t^2 + x^2 + y^2 + z^2$. In fact I will usually work in the Wick-rotated Euclidean space where $x_\mu x^\mu = x^2 + y^2 + z^2 + t^2$.

An analysis in QFT usually follow the following path:

1. The starting point is the classical action, e.g.

$$S = \int d^4x \left(-\frac{1}{2} \partial_\mu \phi \partial^\mu \phi - \frac{1}{2} m^2 \phi^2 \right). \quad (0.1)$$

2. We quantize it into a quantum theory. There are two methods with the same result:

- (a) by the canonical quantization, or
- (b) by the path integral.

3. We now compute the correlation functions

$$\langle \mathcal{O}_1(x_1) \mathcal{O}_2(x_2) \cdots \mathcal{O}_n(x_n) \rangle. \quad (0.2)$$

Again there are various methods:

- (a) perturbative expansions,
- (b) numerical simulations, and
- (c) analytic exact results in some special cases

4. From the correlation functions, we extract experimentally measurable quantities, e.g.

- (a) by Lehmann-Symanzik-Zimmerman (LSZ) reduction to obtain scattering amplitudes,
- (b) by Kubo formula to obtain linear responses.

In this lecture, I will mostly concentrate on the steps 1, 2 and 3, leaving the step 4 for other lecturers. I will mostly use the path integral methods. It is known that there are QFTs which do not have the classical action, i.e. there are no steps 1 and 2. In fact one of my main areas of study is exactly such QFTs, but I do not think I have much time for that in this set of lectures.

1 Recap: path integrals and Feynman diagrams

1.1 Gaussian integration and Wick's theorem

The bedrock of the path integral is the Gaussian integral

$$\int_{-\infty}^{\infty} dx e^{-ax^2/2} = \sqrt{\frac{2\pi}{a}}. \quad (1.1)$$

This is in some sense the 0d quantum field theory: since there are no spacetime directions, the $\partial_\mu \phi \partial^\mu \phi$ term in (0.1) disappears. Then the Gaussian integral (1.1) is the partition function of the 0d QFT. For a function $f(x)$, its expectation value is denoted by

$$\langle f(x) \rangle := \frac{\int_{-\infty}^{\infty} dx f(x) e^{-ax^2/2}}{\int_{-\infty}^{\infty} dx e^{-ax^2/2}}. \quad (1.2)$$

It is clear that $\langle x^{2n+1} \rangle = 0$. We have

$$\langle x^{2n} \rangle = \frac{1}{a^n} (2n-1) \cdot (2n-3) \cdots 3 \cdot 1 \quad (1.3)$$

Note that the right hand side counts the number of possible contractions, weighted by the propagator. This is Wick's theorem. Any other perturbative path integral computations are no more than complicated versions of (1.1) and (1.3).

For example, for an $n \times n$ matrix a_{ij} , we have

$$\int dx_1 \cdots dx_n e^{-a_{ij} x^i x^j / 2} = (2\pi)^{n/2} (\det A)^{-1/2}. \quad (1.4)$$

1.2 Non-convergence of the perturbative expansion

Let us next consider

$$\int_{-\infty}^{\infty} dx e^{-ax^2/2 - bx^4} = \sqrt{\frac{2\pi}{a}} \langle e^{-bx^4} \rangle. \quad (1.5)$$

This integral defines one of the Bessel functions. This can be computed by a perturbative expansion:

$$\langle e^{-bx^4} \rangle = 1 - b \langle x^4 \rangle + \frac{b^2}{2} \langle x^8 \rangle - \frac{b^3}{6} \langle x^{12} \rangle + \cdots = 1 - 3 \frac{b}{a^2} + 105 \frac{b^2}{a^4} - 10395 \frac{b^3}{a^6} + \cdots \quad (1.6)$$

This *does* not converge. One way to see that is to note that the coefficients grow factorially due to (1.3). Therefore the convergence radius is zero. Another way to see that is the following. Suppose the convergence radius is nonzero. Then it would converge also for $b < 0$. But the integral (1.5) clearly does not converge. This is a contradiction.

In perturbative computations of QFT we often encounter infinities. There are mainly two types:

1. The renormalization removes the infinities in the perturbative expansion and produces Taylor series like (1.6) with finite coefficients. In the 0d QFT (1.5) we did not encounter this issue.
2. The resulting perturbative series usually do not converge, essentially due to the same issue as the non-convergence of the perturbative expansion of the Bessel function we saw above. This issue already appears in 0d QFTs, and persists in QFTs in more than 0 dimensions.

We physicists do not often care about this second type of divergence, and evaluate the series by truncating it (usually to the order we are able to compute). Somehow the resulting numerical answer is known to agree with experiments quite well.

1.3 0+1d free scalar field

Let us go on to quantum field theory. The simplest example you learn is the free scalar field, whose action is

$$S = \int d^D x \left(-\frac{1}{2} \partial_\mu \phi \partial^\mu \phi - \frac{1}{2} m^2 \phi^2 \right). \quad (1.7)$$

The associated path integral is

$$Z = \int [D\phi] e^{iS}. \quad (1.8)$$

The Wick-rotated, Euclidean version has the action

$$S_E = \int d^D x \left(\frac{1}{2} \partial_\mu \phi \partial^\mu \phi + \frac{1}{2} m^2 \phi^2 \right). \quad (1.9)$$

and the path integral

$$Z_E = \int [D\phi] e^{-S_E}. \quad (1.10)$$

When $D = 0$, this is really the Gaussian integral we just saw. When $D = 0 + 1$, the action is

$$S = \int dt \left(\frac{1}{2} (\partial_t \phi)^2 - \frac{1}{2} m^2 \phi^2 \right). \quad (1.11)$$

Let us use a different symbol X and ω for what we wrote ϕ and m :

$$S = \int dt \left(\frac{1}{2} (\dot{X})^2 - \frac{1}{2} \omega^2 X^2 \right). \quad (1.12)$$

This is simply the action of a particle of a mass 1 moving along a line X under the potential $V = \omega^2 X^2/2$. Put it differently, this is a quantum mechanical harmonic oscillator.

We all know that the energy eigenvalues are $E = \omega(\frac{1}{2} + n)$, and so the partition function is

$$Z = \text{tr} e^{-\beta H} = e^{-\omega\beta/2} (1 + e^{-\omega\beta} + \dots) = \frac{1}{e^{\omega\beta/2} - e^{-\omega\beta/2}}. \quad (1.13)$$

Let us reproduce it via the path integral. We use the Euclidean version

$$Z \propto \int [DX(t)] \exp\left(-\int dt \left(\frac{1}{2} \partial_t X \partial_t X + \frac{1}{2} \omega^2 X^2 \right)\right) \quad (1.14)$$

where $X(t)$ is now a periodic function on $t \in [0, \beta]$. Following (1.4), this can be schematically written as

$$Z \propto [\det(-\partial_t^2 + \omega^2)]^{-1/2}. \quad (1.15)$$

To actually compute it, we diagonalize the operator inside det by expanding $X(t)$ into Fourier modes $X_n e^{i(2\pi/\beta)nt}$ for $n \in \mathbb{Z}$. We find

$$Z \propto \frac{1}{\omega} \prod_{n \geq 1} \frac{1}{(2\pi/\beta)^2 n^2 + \omega^2} \propto \frac{1}{\beta\omega} \prod_{n \geq 1} \frac{1}{1 + (\frac{\beta\omega}{2\pi n})^2}, \quad (1.16)$$

where we allowed β -dependent but ω -independent infinite proportionality constants. We now note that the infinite product expansion

$$\frac{x}{e^{x/2} - e^{-x/2}} = \prod_{n \geq 1} \frac{1}{1 + (\frac{x}{2\pi n})^2}, \quad (1.17)$$

and see the equivalence between the partition function (1.13) as computed in the canonical quantization and the path integral result of (1.16).

Exercise 1.1. We were not careful about the β -dependence of the infinite proportionality constants. If you would like to be more careful, what should be done?

If we only need the vacuum (or ground state) energy E_0 , we note that $\log Z \sim -\beta E_0$ in the large β limit. Writing $p = 2\pi n/\beta$, we see that there are $\beta dp/(2\pi)$ modes per interval dp , so we have

$$\log Z \sim -\frac{1}{2}\beta \int \frac{dp}{2\pi} \log(p^2 + \omega^2) \sim -\frac{1}{2}\beta \int \frac{dp}{2\pi} \log(1 + \omega^2/p^2) = -\beta \frac{\omega}{2}, \quad (1.18)$$

where \sim means that we added/subtracted β -dependent ω -independent constants.

1.4 ϕ^4 interaction in 0+1d

We now consider the model with the action

$$S = \int dt \left(\frac{1}{2} (\dot{X})^2 - \left(\frac{1}{2} \omega^2 X^2 + \lambda X^4 \right) \right). \quad (1.19)$$

Let us find the ground state energy using the Feynman diagram. We use the Euclidean version:

$$Z = \int [DX] \exp\left(- \int dt (\dot{X})^2 + \frac{1}{2} \omega^2 X^2 + \lambda X^4\right) \quad (1.20)$$

and therefore

$$Z/Z_0 = \langle \exp(- \int dt \lambda X^4) \rangle_0 \quad (1.21)$$

where the subscript $_0$ means that it is evaluated with $\lambda = 0$. We extract the ground state energy E via $\log Z \sim -\beta E$. Writing

$$E = \omega/2 + c_1\lambda + c_2\lambda^2 + \dots, \quad (1.22)$$

we know that

$$Z/Z_0 = -\beta c_1\lambda - (\beta c_2 - \frac{1}{2}(\beta c_1)^2)\lambda^2 - \dots. \quad (1.23)$$

This allows us to read off $c_{1,2,\dots}$ from the computation of Z/Z_0 .

We first Fourier-transform the interaction term to

$$\int dt \lambda X^4 = \int \frac{dp}{2\pi} \frac{dq}{2\pi} \frac{dr}{2\pi} \frac{ds}{2\pi} 2\pi \delta(p+q+r+s) X(p)X(q)X(r)X(s) \quad (1.24)$$

and the propagator is

$$\langle X(p)X(q) \rangle = (2\pi)\delta(p+q) \frac{1}{p^2 + \omega^2}. \quad (1.25)$$

Then the leading contribution to Z/Z_0 is simply given by the Feynman diagram shown in Fig. 1. The value can be evaluated as follows:

$$-\lambda \langle \int dt \lambda X^4 \rangle_0 = -\lambda \int \frac{dp}{2\pi} \frac{dq}{2\pi} \frac{dr}{2\pi} \frac{ds}{2\pi} 2\pi \delta(p+q+r+s) \langle X(p)X(q)X(r)X(s) \rangle_0 \quad (1.26)$$

$$= -3\lambda (2\pi\delta(0)) \underbrace{\left(\int \frac{dp}{2\pi} \frac{1}{p^2 + \omega^2} \right)^2}_{=1/(2\omega)} = -\beta \cdot \frac{3\lambda}{4\omega^2}. \quad (1.27)$$

Here we recall that $2\pi\delta(p) = \int dt e^{itp}$, so by setting $p = 0$ we know that $2\pi\delta(p)$ acts as $\beta = \int dt$ in the large β limit. Comparing with (1.22) and (1.23), we find that the ground state energy is

$$E = \frac{\omega}{2} + \frac{3\lambda}{4\omega^2} \quad (1.28)$$

to this order.

In the next order, we see three Feynman diagrams, see Fig. 1 again. The diagram a) contributes exactly as $1/2$ the square of (1.27), so from (1.23) we see that the second order contribution to E has two contribution, one is

$$\text{b)} = \beta\lambda^2 \cdot \frac{1}{2} \cdot 6 \cdot 6 \left(\int \frac{dp}{2\pi} \frac{1}{p^2 + \omega^2} \right)^2 \underbrace{\left(\int \frac{dp}{2\pi} \frac{1}{(p^2 + \omega^2)^2} \right)}_{1/(4\omega^3)} = \beta \frac{9\lambda^2}{4\omega^5} \quad (1.29)$$

and another is

$$\text{c)} = \beta \cdot \lambda^2 \frac{1}{2} \cdot 24 \underbrace{\int \frac{dpdqdr}{(2\pi)^3} \frac{1}{p^2 + m^2} \frac{1}{q^2 + m^2} \frac{1}{r^2 + m^2} \frac{1}{(p+q+r)^2 + m^2}}_{=1/(32\omega^5)} = \beta \frac{3\lambda^2}{8\omega^5}. \quad (1.30)$$

In total, we find

$$E = \frac{\omega}{2} + \frac{3\lambda}{4\omega^2} - \frac{21\lambda^2}{8\omega^5} + \dots \quad (1.31)$$

to this order. Reassuringly, what we obtained so far reproduces the results obtained from the standard perturbation theory in quantum mechanics, see Sec. 1.5 and in particular (1.47).

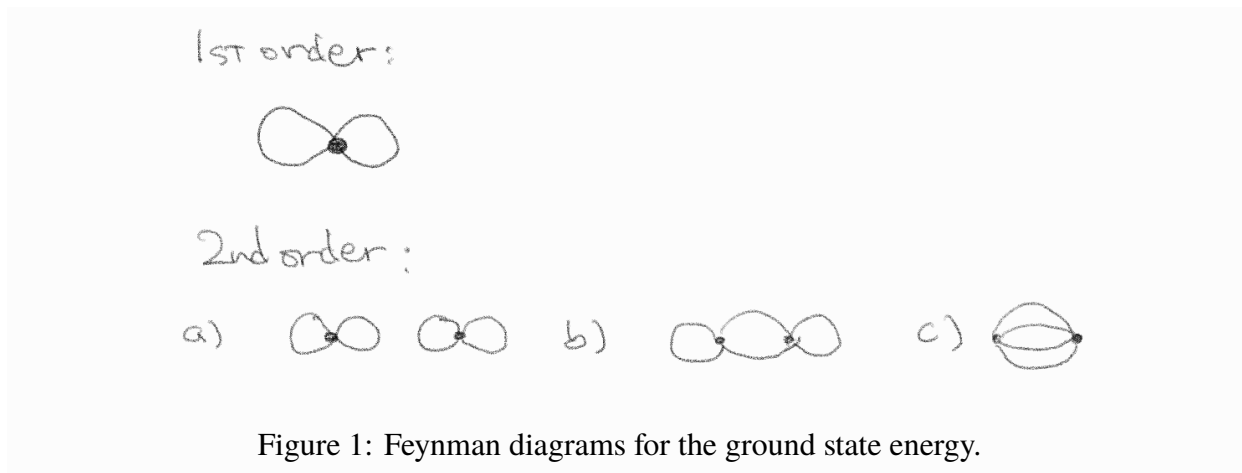


Figure 1: Feynman diagrams for the ground state energy.

Exercise 1.2. Carry out the computation in the next order, i.e. to the order λ^3 .

Exercise 1.3. Write a program in your favorite computer algebra system to perform the computation in an arbitrary order, and compare its result against the result from the standard perturbation theory in quantum mechanics to your satisfaction.

1.5 Aside: Brillouin-Wigner perturbation theory in quantum mechanics

Here we quickly recall Brillouin-Wigner perturbation theory in quantum mechanics. (I do not plan to talk about this subsection in the actual lecture.)

We would like to solve

$$(H_0 + gV)|\psi_n(g)\rangle = E_n(g)|\psi_n(g)\rangle \tag{1.32}$$

where $\psi_n(g)$ is the n -th eigenstate, with the normalization

$$\langle\psi_n|\psi_n(g)\rangle = 1. \tag{1.33}$$

We start by rewriting (1.32) as

$$(E_n(g) - H_0)|\psi_n(g)\rangle = gV|\psi_n(g)\rangle. \tag{1.34}$$

We first hit both sides of (1.34) from the right by $\langle\psi_n|$ and find

$$E_n(g) - E_n = \langle\psi_n|gV|\psi_n(g)\rangle. \tag{1.35}$$

This means that if we know $\psi_n(g)$ up to g^i , we know $E_n(g)$ up to g^{i+1} .

We next apply $\langle \psi_m |$ ($m \neq n$):

$$(E_n(g) - E_m) \langle \psi_m | \psi_n(g) \rangle = g \langle \psi_m | V | \psi_n(g) \rangle \quad (1.36)$$

which means

$$|\psi_n(g)\rangle = |\psi_n\rangle + g \left(\sum_{m \neq n} |\psi_m\rangle \frac{1}{E_n(g) - E_m} \langle \psi_m | \right) V |\psi_n(g)\rangle. \quad (1.37)$$

We now introduce

$$R_n(g) = \sum_{m \neq n} |\psi_m\rangle \frac{1}{E_n(g) - E_m} \langle \psi_m | \quad (1.38)$$

with which we can write

$$|\psi_n(g)\rangle = |\psi_n\rangle + g R_n(g) V |\psi_n(g)\rangle \quad (1.39)$$

$$= |\psi_n\rangle + g R_n(g) V |\psi_n\rangle + g^2 R_n(g) V R_n(g) V |\psi_n(g)\rangle \quad (1.40)$$

$$= |\psi_n\rangle + g R_n(g) V |\psi_n\rangle + g^2 R_n(g) V R_n(g) V R_n(g) V |\psi_n(g)\rangle \quad (1.41)$$

ad infinitum. $E_n(g)$ can then be found by (1.35).

So far we did not make any approximation. We perform the same expansion repeatedly and replace the final $|\psi_n(g)\rangle$ by $|\psi_n\rangle$. We obtain

$$|\psi_n(g)\rangle = |\psi_n\rangle + g R_n(g) V |\psi_n\rangle + g^2 R_n(g) V R_n(g) V R_n(g) V |\psi_n\rangle + \dots \quad (1.42)$$

We now multiply $\langle \psi_n | gV$ from the right and obtain

$$E_n(g) = E_n + g \langle \psi_n | V | \psi_n \rangle + g^2 \langle \psi_n | V R_n(g) V | \psi_n \rangle + g^3 \langle \psi_n | V R_n(g) V R_n(g) V | \psi_n \rangle + \dots \quad (1.43)$$

We note that the denominator of $R_n(g)$ contains $E_n(g)$ itself, which needs to be expanded using (1.43). This is the Brillouin-Wigner perturbation, which is more economical than the very basic Rayleigh-Schrödinger perturbation.

For example, the g^3 -contribution to $E_n(g)$ is a sum of $g^3 \langle \psi_n | V R_n(g) V R_n(g) V | \psi_n \rangle$ where $R_n(g)$ is replaced by $R_n(0)$ and $g^2 \langle \psi_n | V R_n(g) V | \psi_n \rangle$ where we evaluate $R_n(g)$ to the first order in g . The latter is given by

$$\frac{1}{E_n(g) - E_m} = \frac{1}{E_n - E_m} - \frac{1}{E_n - E_m} (E_n(g) - E_n) \frac{1}{E_n - E_m} \quad (1.44)$$

$$= \frac{1}{E_n - E_m} - \frac{1}{E_n - E_m} g \langle \psi_n | V | \psi_n \rangle \frac{1}{E_n - E_m} + O(g^2). \quad (1.45)$$

In total, we find

$$E_n^{(3)} = \langle \psi_n | V R_n(0) V R_n(0) V | \psi_n \rangle - \langle \psi_n | V R_n(0) R_n(0) V | \psi_n \rangle \langle \psi_n | V | \psi_n \rangle. \quad (1.46)$$

For an actual computation it is useful to use any computer algebra system of your liking.

Exercise 1.4. Perform this computation for the example where $H = \frac{1}{2}(p^2 + x^2) + gx^4$.

Answer. My sample implementation follows:

```
(* tell Mathematica how operators act on the states*)
Act[{c___, a_ + b_}, psi_] := Act[{c, a}, psi] + Act[{c, b}, psi];
Act[a_, psi_ + phi_] := Act[a, psi] + Act[a, phi];
(* Ket[n] is the unperturbed n-th eigenstate*)
Act[a_, c_ Ket[n_]] := c Act[a, Ket[n]];
(* A is the annihilator and AD is the creator*)
Act[{c___, A}, Ket[n_]] := Sqrt[n] Act[{c}, Ket[n - 1]];
Act[{c___, AD}, Ket[n_]] := Sqrt[n + 1] Act[{c}, Ket[n + 1]];
Act[{}, psi_] := psi;
(* This is the X^4 operator *)
V[psi_] := Act[{A + AD, A + AD, A + AD, A + AD}, psi]/4;
(* Define the operation to extract the coefficient of |0>*)
Proj[psi_ + phi_] := Proj[psi] + Proj[phi];
Proj[c_ Ket[n_]] := c Proj[Ket[n]];
Proj[Ket[n_]] := If[n == 0, 1, 0];
(* Define the operator R*)
R[order_, psi_ + phi_] := R[order, psi] + R[order, phi];
R[order_, c_ Ket[n_]] := c R[order, Ket[n]];
R[order_, Ket[n_]] := If[n == 0, 0, Ket[n]/(Energy[order] - (n + 1/2))];
(* State[n] is the perturbed n-th eigenstate *)
(* Energy[n] is the perturbed n-th eigenvalue *)
State[0] = Ket[0];
Energy[0] = 1/2;
VState[order_] := VState[order] =
  V[State[order] // Normal // Expand] // Expand;
Energy[order_] := Energy[order] =
  Energy[0] + g Proj[VState[order - 1]] + O[g]^(order + 1) // Normal;
State[order_] := State[order] =
  State[0] + g R[order - 1, VState[order - 1]] + O[g]^(order + 1) // Normal;
(*This will compute the perturbative series to g^20*)
Energy[20]
```

$$E_0(g) = \frac{1}{2} + \frac{3g}{4} - \frac{21g^2}{8} + \frac{333g^3}{16} - \frac{30885g^4}{128} + \frac{916731g^5}{256} - \frac{65518401g^6}{1024} + \frac{2723294673g^7}{2048} - \frac{1030495099053g^8}{32768} + \frac{54626982511455g^9}{65536} - \frac{6417007431590595g^{10}}{262144} + O(g^{11}). \quad (1.47)$$

1.6 Vacuum energy of a 3+1d free scalar field

Let us finally talk about 3+1d QFT. We simply consider the free massive scalar field with the action

$$S = \int d^4x \left(-\frac{1}{2} \partial_\mu \phi \partial^\mu \phi - \frac{1}{2} m^2 \phi^2 \right). \quad (1.48)$$

The Euclidean path integral is then

$$Z \propto \int [D\phi] \exp\left(- \int d^4x \left(\frac{1}{2} |\partial_\mu \phi|^2 + \frac{1}{2} m^2 \phi^2 \right)\right). \quad (1.49)$$

We can evaluate $\log Z$ as in 0+1d.

Here we put the system on a large spacetime box of size $L_{1,2,3,4}$. Then we see it behaves as $\log Z \sim L_1 L_2 L_3 L_4 E$ where E is the energy density per spatial volume, where

$$E \sim \frac{1}{2} \int \frac{d^4p}{(2\pi)^4} \log(p^2 + m^2) \sim \frac{1}{2} \int \frac{d^4p}{(2\pi)^4} \log(1 + m^2/p^2). \quad (1.50)$$

where \sim allows an addition/subtraction of m -independent constants, see (1.18) for the 0+1d case.

We see that this still diverges. Taking the derivative w.r.t. m^2 three times, we find that the integral converges and

$$\left(\frac{\partial}{\partial m^2}\right)^3 E = \frac{1}{2} \int \frac{d^4p}{(2\pi)^4} \frac{2}{(p^2 + m^2)^3} = \frac{2\pi^2}{(2\pi)^2} \int_0^\infty \frac{p^3 dp}{(p^2 + m^2)^3} = \frac{1}{32\pi^2 m^2}. \quad (1.51)$$

Integrating w.r.t. m^2 three times, this means that the vacuum energy density is

$$E = am^4 + \Lambda^2 m^2 + \frac{1}{64\pi^2} m^4 \log \frac{m^2}{\Lambda'^2}. \quad (1.52)$$

where a , Λ and Λ' are three integration constants (where a is dimensionless and Λ , Λ' has a dimension of energy.)

Note that this is in principle measurable, since the vacuum energy density acts as cosmological constants. Therefore, if the mass m slowly varies in a cosmological timescale, this would manifest as a slow variation of the cosmological constant. The same computation also appears as an intermediate step in the evaluation of the Coleman-Weinberg potential, which describes how the scalar potential is quantum mechanically modified.

Exercise 1.5. Perform the computation in the different spacetime dimensions. What is the physical interpretation in the very low dimensional case, say in 0d or in 0+1d?

1.7 The Casimir effect

Let us next consider a massless scalar field in 0+1d, but with the constraint that the spatial direction is periodically identified, $x \sim x + L$. The action is as always

$$S = \int dt dx \left(-\partial_\mu \phi \partial^\mu \phi - V_0 \right) \quad (1.53)$$

where we included a constant term V_0 . The vacuum energy V is then

$$V = LV_0 + \sum_{n=1}^{\infty} E_n, \quad E_n = \frac{2\pi n}{L} \quad (1.54)$$

We regularize by writing it as

$$V = LV_0 + \sum_{n=1}^{\infty} e^{-E_n/\Lambda} E_n = LV_0 + \frac{L\Lambda^2}{2\pi} - \frac{1}{12} \frac{2\pi}{L} + O(\Lambda^{-4}), \quad (1.55)$$

where Λ is a large energy cutoff. This means that we can renormalize it by taking $V_0 = -\Lambda^2/(2\pi)$, and the vacuum energy is found to be

$$V = -\frac{1}{12} \frac{2\pi}{L}. \quad (1.56)$$

Exercise 1.6. Perform the same computation in the case of a massive scalar field.

The same computation can be performed for a massless free scalar in 3+1d, with one spatial direction restricted to $0 < x < L$, with the Dirichlet boundary condition $\phi(x=0) = \phi(x=L) = 0$. The energy density turns out to be

$$E = -\frac{\pi^2}{1440L^4}. \quad (1.57)$$

Exercise 1.7. Compute it.

For the electromagnetic field within two perfectly conducting plates, we need to multiply it by a factor of two, accounting for two polarizations. This is the Casimir effect, and has been measured, see Fig. 2.

We now know that we have a nonzero cosmological constant. The measured value is about 10^{-29}g/cm^3 . The Casimir energy of the electromagnetic field will be comparable to this when

$$\frac{\hbar c \pi^2}{720L^4} \sim 10^{-29} \text{g}c^2/\text{cm}^3, \quad (1.58)$$

which is when

$$L \sim 25 \mu\text{m}. \quad (1.59)$$

Exercise 1.8. Is my computation correct? If so, why is the current cosmological constant, which is related to the largest scale in the universe, has such an ordinary looking value as (1.59)?

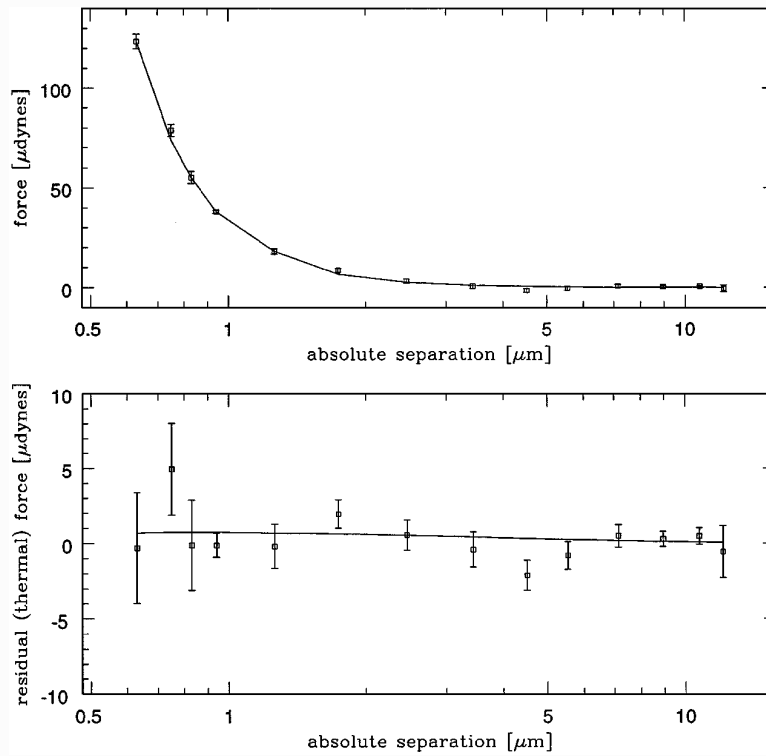


FIG. 4. Top: All data with electric force subtracted, averaged into bins (of varying width), compared to the expected Casimir force for a 11.3 cm spherical plate. Bottom: Theoretical Casimir force, without the thermal correction, subtracted from top plot; the solid line shows the expected residuals.

Figure 2: Measurement of the Casimir force by Lamoreaux, taken from [Lam97].

2 The Lagrangian of the Standard Model

2.1 Gauge fields

Let $A_\mu(x)$ be a vector field ($\mu = 0, 1, 2, 3$) whose components are $N \times N$ Hermitean matrices. This means, for example, $A_{\mu=0}(x)$ is an $N \times N$ Hermitean matrix with elements $(A_{\mu=0})_j^i(x)$, $i = 1, \dots, N$, and each element $(A_\mu)_j^i(x)$ is a function on the spacetime.

We consider the combination

$$D_\mu := \partial_\mu + iA_\mu \quad (2.1)$$

(where the first term ∂_μ is implicitly multiplied with the identity matrix). We now consider $g(x)$ which assigns $N \times N$ unitary matrices at each point x in the spacetime. We then define the gauge transformation by $g(x)$

$${}^g D_\mu := g D_\mu g^{-1}. \quad (2.2)$$

We define ${}^g A_\mu$ by the formula

$${}^g D_\mu = \partial_\mu + i{}^g A_\mu \quad (2.3)$$

and find

$${}^g A_\mu = g A_\mu g^{-1} + i g \partial_\mu g^{-1}. \quad (2.4)$$

This is the $U(N)$ gauge field.

Exercise 2.1. Check that ${}^g A_\mu$ is Hermitean.

We can impose the conditions that $\text{tr } A_\mu = (A_\mu)_i^i = 0$. Accordingly, we restrict $g(x)$ so that $\det g = 1$; such matrices are called special unitary. This guarantees that $\text{tr } {}^g A_\mu = 0$. This is the $SU(N)$ gauge field.

Exercise 2.2. Check that $\text{tr } {}^g A_\mu = 0$.

We now define

$$iF_{\mu\nu} := [D_\mu, D_\nu] = i(\partial_\mu A_\nu - \partial_\nu A_\mu + i[A_\mu, A_\nu]). \quad (2.5)$$

We can easily compute its gauge variation to be

$${}^g F_{\mu\nu} = g F_{\mu\nu} g^{-1} \quad (2.6)$$

and therefore

$$\text{tr } F_{\mu\nu} F^{\mu\nu}, \quad \text{tr } F_{\mu\nu} F_{\rho\sigma} \epsilon^{\mu\nu\rho\sigma} \quad (2.7)$$

are both gauge invariant.

The classical action of the $SU(N)$ gauge field is then taken to be

$$S_{\text{Yang-Mills}} = - \int d^4x \left[\frac{1}{2g^2} \text{tr } F_{\mu\nu} F^{\mu\nu} + \frac{\theta}{32\pi^2} \text{tr } F_{\mu\nu} F_{\rho\sigma} \epsilon^{\mu\nu\rho\sigma} \right]. \quad (2.8)$$

g is the coupling constant and θ is the theta angle.

- We will see later that due to quantum effects g changes as we change the scale we measure the system.
- We will also see later that θ and $\theta + 2\pi$ cannot be distinguished, and therefore θ is an angular parameter.
- Note that this is the *math* normalization. Physicists usually define $A_\mu^{\text{math}} = gA_\mu^{\text{phys}}$ so that g does not appear as the coefficient of the $\text{tr } F_{\mu\nu} F^{\mu\nu}$ term. The math normalization obscures some physics but makes various formulas somewhat shorter.

A $U(1)$ gauge field is simply the Maxwell field. Indeed, for $N = 1$, an $N \times N$ matrix is a number, and $g(x)$ can be written as $g(x) = e^{i\chi(x)}$. Then the transformation law (2.4) becomes

$${}^g A_\mu = A_\mu + i\partial_\mu \chi. \quad (2.9)$$

The commutator term in (2.5) also drops out, and we simply have

$$F_{\mu\nu} = \partial_\mu A_\nu - \partial_\nu A_\mu. \quad (2.10)$$

The Lagrangian is usually written as

$$S_{\text{Maxwell}} = - \int d^4x \left[\frac{1}{4e^2} F_{\mu\nu} F^{\mu\nu} + \frac{\theta}{32\pi^2} F_{\mu\nu} F_{\rho\sigma} \epsilon^{\mu\nu\rho\sigma} \right] \quad (2.11)$$

where e is the electric charge and θ is again the theta angle.

- The difference in the coefficient $1/2$ in (2.8) and $1/4$ in (2.11) are not typos but follows the standard convention in the literature.
- θ is again an angular variable.
- We will see later that θ affects the electric charge of a magnetic monopole, so it is in principle a measurable quantity. It is also known that within a time-reversal-invariant topological insulator, we have $\theta_{\text{inside}} - \theta_{\text{outside}} = \pi$.

Here we only considered the case when the gauge transformation $g(x)$ is a (special) unitary $N \times N$ matrix, thus forming the group $U(N)$ or $SU(N)$. More generally, we can use whatever compact Lie group G and its Lie algebra \mathfrak{g} , so that the components $A_{\mu=0,1,2,3}$ of the gauge fields now take values in \mathfrak{g} . The formulations above can be readily adapted to this more general setting, after taking care of the normalization of tr .

2.2 Why gauge invariance?

Gauge invariance is required for a massless vector field. Let us briefly recall why. The logic here is not very precise. See textbooks for details. (The keywords are the BRST transformation and the Kugo-Ojima condition.)

In a (naive) covariant quantization, we introduce creation operators $a_\mu^\dagger(k)$ for each μ . The inner product is

$$\langle 0|a_\mu a_\nu^\dagger|0\rangle \propto \eta_{\mu\nu}. \quad (2.12)$$

In particular, when the norm of the spatial component $\langle 0|a_{\mu=3}a_{\nu=3}^\dagger|0\rangle$ is positive, the norm of the temporal component $\langle 0|a_{\mu=0}a_{\nu=0}^\dagger|0\rangle$ is negative. This is bad, because the norm is related to the probability, and should be positive.

In a gauge-invariant system, one way out is to declare that any modes which arise from gauge transformation are decoupled and can be removed. For simplicity, say $k_\mu = (k, 0, 0, k)$. The equation of motion $\partial^\mu F_{\mu\nu} = 0$ implies that $k^\mu a_\mu^\dagger = 0$, which means that only a_1^\dagger , a_2^\dagger , and $a_0^\dagger + a_3^\dagger$ are kept.

Then, the gauge transformation $\partial_\mu \chi$ can create the linear combination $a_0^\dagger + a_{\mu=3}^\dagger$. In the end, what remains is $a_{\mu=1}^\dagger|0\rangle$ and $a_{\mu=2}^\dagger|0\rangle$, which are the two modes transverse to k^μ and describe the two transverse polarizations. Effectively, one mode $a_0^\dagger|0\rangle$ with a negative norm and another mode $a_3^\dagger|0\rangle$ with a positive norm with exactly the same magnitude paired up and cancelled.

With gauge invariance, a detailed version of this analysis can be done in the presence of interactions, and can show that the norm of the Hilbert space is positive definite, and the probability is positive. Without gauge invariance, we eventually run into a problem of negative probability.

2.3 Coupling gauge fields to scalars

Consider N complex scalar fields $\phi^i(x)$ ($i = 1, \dots, N$). The gauge transformation $g_j^i(x)$ acts as by a matrix multiplication:

$${}^g\phi(x) := g\phi(x). \quad (2.13)$$

This is called a scalar in the fundamental representation of $U(N)$ or $SU(N)$.

From the definition of the covariant derivative, we see

$${}^g(D_\mu\phi(x)) = gD_\mu\phi(x). \quad (2.14)$$

Then the action

$$S_{\text{scalar kinetic term}} = - \int d^4x (D_\mu\phi(x))^\dagger D^\mu\phi(x) \quad (2.15)$$

is gauge invariant, since $g^\dagger g = \text{id}$ for a unitary matrix.

We can also add a gauge-invariant potential $V(\phi)$ to this action. One example is

$$V(\phi) = m^2\phi^\dagger\phi + \frac{1}{2}\lambda^2(\phi^\dagger\phi)^2; \quad (2.16)$$

then the total action is

$$S_{\text{scalar}} = - \int d^4x [(D_\mu\phi(x))^\dagger D^\mu\phi(x) + V(\phi(x))] \quad (2.17)$$

More generally, consider a general gauge group G and a representation ρ which is n -dimensional. This means that $\rho(g)$ is an $n \times n$ matrix and $\rho(g)\rho(h) = \rho(gh)$. We introduce $\phi_a(x)$ where $a = 1, \dots, n$. The gauge transformation is then defined by

$${}^g\phi(x) := \rho(g)\phi(x). \quad (2.18)$$

This is called a complex scalar in the representation ρ .

The covariant derivative we should use is

$$\rho(D_\mu) := \partial_\mu + i\rho(A_\mu) \quad (2.19)$$

where the representation matrix $\rho(X)$ for a Lie algebra generator $X \in \mathfrak{g}$ satisfies

$$\rho(\underbrace{e^{i\epsilon X}}_{=g}) = \text{id} + i\epsilon\rho(X) + \dots \quad (2.20)$$

This guarantees that

$$\rho({}^gD_\mu) = \rho(g)\rho(D_\mu)\rho(g)^{-1} \quad (2.21)$$

and that the action

$$S_{\text{scalar kinetic term}} = - \int d^4x (\rho(D_\mu)\phi(x))^\dagger \rho(D^\mu)\phi(x) \quad (2.22)$$

is gauge invariant. Since we cannot use anything other than $\rho(D_\mu)$ in the expression above, we do not typically write ρ here, and simply use (2.15) instead.

Consider in particular the case $G = U(1)$, and the 1-dimensional representation $\rho_q(g) = g^q$; note g is a 1×1 unitary matrix, i.e. a complex number with absolute value 1. It is easy to check that $\rho_q(gh) = \rho_q(g)\rho_q(h)$. Consider a scalar ϕ in this representation; such a field transforms as

$${}^g\phi(x) = g(x)^q\phi(x) \quad (2.23)$$

under the gauge transformation. What should be the covariant derivative? From (2.20) we see $\rho_q(X) = qX$ for a Lie algebra generator. This means that the covariant derivative in this case is

$$\rho_q(D_\mu)\phi = (\partial_\mu + iqA_\mu)\phi. \quad (2.24)$$

Thus we see that this number q is the electric charge of the field ϕ .

Next, consider the case $G = U(1) \times SU(2)$, whose element is parameterized by $g = (g_1, g_2)$ where g_1 is a complex number of absolute value 1 and g_2 is a 2×2 special unitary matrix. We can introduce a two-component scalar field with the transformation law

$${}^g\phi(x) = \underbrace{(g_1)^q}_{=\rho(g)} g_2 \phi(x). \quad (2.25)$$

This is a field of $U(1)$ charge q in the doublet representation of $SU(2)$. The corresponding covariant derivative has the form

$$\rho(D_\mu)\phi(x) = (\partial_\mu + iqA_\mu^{U(1)} + iA_\mu^{SU(2)})\phi(x). \quad (2.26)$$

Note that the first two terms on the right hand side are implicitly multiplied by the 2×2 identity matrix, and $A_\mu^{SU(2)}$ are themselves 2×2 matrices.

2.4 Two-component spinors

Let us consider the group $\text{SL}(2, \mathbb{C})$ of 2×2 complex matrices g of determinant 1.¹ This acts naturally on Hermitean 2×2 matrices X by the formula

$$X \mapsto {}^g X = gXg^\dagger, \quad (2.27)$$

since gXg^\dagger is also Hermitean.

We parameterize the Hermitean matrix X as

$$X = \begin{pmatrix} t + z & x + iy \\ x - iy & t - z \end{pmatrix} \quad (2.28)$$

We note that

$$\det X = t^2 - x^2 - y^2 - z^2 \quad (2.29)$$

is the minus of the standard norm of the Minkowski space $\mathbb{R}^{3,1}$. Therefore

$$\det({}^g X) = \det(gXg^\dagger) = (\det g)(\det X)(\det g^\dagger) = \det X. \quad (2.30)$$

We define $({}^g x)^\mu$ in terms of ${}^g X$ by the formula (2.28). This then means that

$$({}^g x)^\mu ({}^g x)_\mu = x^\mu x_\mu, \quad (2.31)$$

which should be a Lorentz transformation

$$x^\mu \mapsto ({}^g x)^\mu := \Lambda(g)^\mu_\nu x^\nu. \quad (2.32)$$

Note that $g = -1 \in \text{SL}(2, \mathbb{C})$ keeps X intact. Therefore we constructed a homomorphism $\text{SL}(2, \mathbb{C}) \rightarrow \text{SO}(3, 1)$, which is 2 to 1. For example, consider

$$\text{SL}(2, \mathbb{C}) \ni g = \begin{pmatrix} e^{i\theta} & 0 \\ 0 & e^{-i\theta} \end{pmatrix} \quad (2.33)$$

By computing gXg^{-1} , one finds that this corresponds to a 2θ rotation around the z axis. So a 2π rotation around the z axis corresponds to $-1 \in \text{SL}(2, \mathbb{C})$ and the 4π rotation around the z axis finally comes back to $1 \in \text{SL}(2, \mathbb{C})$.

Exercise 2.3. What does the element

$$\text{SL}(2, \mathbb{C}) \ni g = \begin{pmatrix} e^\beta & 0 \\ 0 & e^{-\beta} \end{pmatrix} \quad (2.34)$$

do as a Lorentz transformation?

¹In this subsection, g is for the spacetime transformation. g is a gauge transformation in most of the other parts of the notes.

Exercise 2.4. Write down the Lorentz transformation $\Lambda(g)$ in (2.32) explicitly.

Two-component spinors are two-dimensional representations of this $SL(2, \mathbb{C})$.² A left-handed spinor ψ^α ($\alpha = 1, 2$) is transformed as

$${}^g\psi^\alpha = g^\alpha_\beta \psi^\beta \quad (2.35)$$

and its complex conjugate, the right-handed spinor $\bar{\psi}^{\dot{\alpha}}$ ($\dot{\alpha} = \dot{1}, \dot{2}$) is transformed as

$${}^g\bar{\psi}^{\dot{\alpha}} = \bar{g}^{\dot{\alpha}}_\beta \bar{\psi}^{\dot{\beta}}. \quad (2.36)$$

Exercise 2.5. Is my convention correct? The choice of $\pm i$ in front of y in (2.28) ties the handedness to the discussion, so it should be possible to check if ψ^α as I defined here leads to massless left-handed fermion. I confess I have not checked this.

Here we adopt the rule where an index α or $\dot{\alpha}$ gains or loses the dot when it is taken outside the bar for the complex conjugation: e.g.

$$\overline{g^\alpha_\beta} = \bar{g}^{\dot{\alpha}}_{\dot{\beta}}. \quad (2.37)$$

In my convention

$$\overline{\bar{\psi}^{\dot{\alpha}}} = \bar{\psi}^\alpha. \quad (2.38)$$

For two spinors ψ^α and χ^α , let us consider the combination

$$\epsilon_{\alpha\beta} \psi^\alpha \chi^\beta := \det(\psi, \chi) := \det \begin{pmatrix} \psi^1 & \chi^1 \\ \psi^2 & \chi^2 \end{pmatrix}. \quad (2.39)$$

Here

$$\epsilon_{11} = \epsilon_{22} = 0, \quad \epsilon_{12} = -\epsilon_{21} = 1. \quad (2.40)$$

This is invariant under the Lorentz transformation, since

$$\epsilon_{\alpha\beta} ({}^g\psi)^\alpha ({}^g\chi)^\beta = \det(g\psi, g\chi) = \det(g \cdot (\psi, \chi)) = \det g \det(\psi, \chi) = \epsilon_{\alpha\beta} \psi^\alpha \chi^\beta. \quad (2.41)$$

We raise and lower the spinor indices using this $\epsilon_{\alpha\beta}$, just as we raise and lower the vector indices μ using $\eta_{\mu\nu}$. (One needs to be extra careful about the signs because ϵ is antisymmetric while η was symmetric. I will not be precise about the signs below.) Then we have

$$\epsilon_{\alpha\beta} \psi^\alpha \chi^\beta = \psi^\alpha \chi_\alpha. \quad (2.42)$$

In this notation, the transformation (2.27) becomes

$$X^{\alpha\dot{\alpha}} \mapsto ({}^gX)^{\alpha\dot{\alpha}} = g^\alpha_\beta \bar{g}^{\dot{\alpha}}_{\dot{\beta}} X^{\beta\dot{\beta}}. \quad (2.43)$$

²I would often call them Weyl fermions but particle phenomenologists would probably disagree with my usage.

We also write (2.28) as

$$X^{\alpha\dot{\alpha}} = \sigma_{\mu}^{\alpha\dot{\alpha}} x^{\mu} \quad (2.44)$$

where

$$\sigma_{\mu=0} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \sigma_{\mu=1} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \sigma_{\mu=2} = \begin{pmatrix} 0 & i \\ -i & 0 \end{pmatrix}, \quad \sigma_{\mu=3} = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}. \quad (2.45)$$

We can now write down the Lagrangian for a massive two-component fermion ψ^{α} :

$$S_{\text{Weyl}} = \int d^4x [i\bar{\psi}^{\dot{\alpha}} \sigma_{\alpha\dot{\alpha}}^{\mu} \partial_{\mu} \psi^{\alpha} + \frac{m}{2} \psi^{\alpha} \psi_{\alpha} + \frac{\bar{m}}{2} \bar{\psi}_{\dot{\alpha}} \bar{\psi}^{\dot{\alpha}}]. \quad (2.46)$$

Note that $m\psi^{\alpha}\psi_{\alpha}$ is nonzero because the fermions ψ anticommute:

$$\psi^{\alpha}\psi_{\alpha} = \epsilon_{\alpha\beta}\psi^{\alpha}\psi^{\beta} = \psi^1\psi^2 - \psi^2\psi^1 = 2\psi^1\psi^2. \quad (2.47)$$

This type of the mass term involving a single two-component spinor is called a Majorana mass term, and a massive fermion of this type is called a Majorana fermion.

A Dirac fermion Ψ has four components, and decomposes into two two-component spinors:

$$\Psi = \begin{pmatrix} \psi_{\alpha} \\ \bar{\chi}^{\dot{\alpha}} \end{pmatrix} \quad (2.48)$$

and the gamma matrices in this basis are

$$\gamma^{\mu} = \begin{pmatrix} 0 & \sigma^{\mu\dot{\alpha}\alpha} \\ \sigma_{\alpha\dot{\alpha}}^{\mu} & 0. \end{pmatrix} \quad (2.49)$$

Then the Dirac action can be decomposed as follows:

$$S_{\text{Dirac}} = \int d^4x [i\bar{\Psi}\gamma^{\mu}\partial_{\mu}\Psi - m\bar{\Psi}\Psi] \quad (2.50)$$

$$= \int d^4x [i\bar{\psi}^{\dot{\alpha}}\sigma_{\alpha\dot{\alpha}}^{\mu}\partial_{\mu}\psi^{\alpha} + i\bar{\chi}^{\dot{\alpha}}\sigma_{\alpha\dot{\alpha}}^{\mu}\partial_{\mu}\chi^{\alpha} - m\chi^{\alpha}\psi_{\alpha} - m\bar{\psi}_{\dot{\alpha}}\bar{\chi}^{\dot{\alpha}}]. \quad (2.51)$$

The mass term now connects two two-component spinors ψ and χ .

In the following we often omit the indices when there are only one sensible way of putting them. This needs practice but simplifies writing.

Exercise 2.6. Check this decomposition. (I am sorry that the signs and factors written above are most probably wrong.)

2.5 Coupling gauge fields to fermions

It is not really different from the coupling of gauge fields to scalars. Let us say that a gauge transformation $g(x)$ acts on a fermion $\psi^{i\alpha}(x)$ via a unitary representation ρ as

$${}^g\psi^{i\alpha} = \rho(g)_j^i \psi^{j\alpha}. \quad (2.52)$$

Then

$$\rho({}^g D_\mu) {}^g\psi = \rho(g)_j^i (D_\mu \psi)^j. \quad (2.53)$$

Combining with

$${}^g\bar{\psi}^i = \overline{\rho(g)_j^i \psi^j}, \quad (2.54)$$

one finds that

$$\int d^4x [i\bar{\psi}\sigma D_\mu\psi] \quad (2.55)$$

is a gauge-invariant Lagrangian.

The mass term will be of the form

$$m c_{ij} \epsilon_{\alpha\beta} \psi^{i\alpha} \psi^{j\beta}. \quad (2.56)$$

Since $\epsilon_{\alpha\beta}$ is antisymmetric and ψ anticommutes, c_{ij} needs to be symmetric. c_{ij} also needs to be gauge-invariant. Therefore, the mass term is possible only when the representation $\rho(g)_j^i$ admits such an invariant symmetric two-index tensor. This is equivalent to the condition that $\rho(g)_j^i$ is a strictly real representation.

For example, let $G = U(1)$ and consider a two-component fermion ψ of charge $q \neq 0$:

$${}^g\psi = g^q \psi \quad (2.57)$$

where $g \in U(1)$ is a complex number of absolute value one. We then have

$${}^g(\psi^\alpha \psi_\alpha) = g^{2q} \psi^\alpha \psi_\alpha, \quad (2.58)$$

and the mass term is not gauge invariant, and therefore is not allowed. To write down a gauge-invariant mass term, one introduces χ of charge $-q$ with the transformation

$${}^g\chi = g^{-q} \chi. \quad (2.59)$$

Then

$${}^g(\psi^\alpha \chi_\alpha) = \psi^\alpha \chi_\alpha \quad (2.60)$$

is gauge invariant. If we start from a Dirac spinor Ψ and assigns charge $+q$

$${}^g\Psi = g^q \Psi, \quad (2.61)$$

we obtain a pair of two component spinors ψ, χ of charge $+q, -q$ via the decomposition (2.48).

Exercise 2.7. Consider a fermion $\psi^{i\alpha}$ where $i = 1, 2$ is the gauge index for $SU(2)$, i.e. it has the gauge transformation

$${}^g\psi^i = g_j^i \psi^j \quad (2.62)$$

for 2×2 special unitary matrices g . Can you write down a mass term?

2.6 Coupling scalars to fermions

It is also possible to couple scalars to the fermions in a gauge-invariant manner. For example, let us go back to (2.58). We can compensate this transformation if we have a scalar field ϕ of charge $-2q$, with the transformation

$${}^g\phi = g^{-2q}\phi. \quad (2.63)$$

Indeed, the term $\phi\psi^\alpha\psi_\alpha$ is now gauge invariant

$${}^g(\phi\psi^\alpha\psi_\alpha) = \phi\psi^\alpha\psi_\alpha, \quad (2.64)$$

and therefore we can consider the gauge-invariant interaction term

$$S_{\text{Yukawa}} = \int d^4x [y\phi\psi^\alpha\psi_\alpha + \bar{y}\bar{\phi}\bar{\psi}_\alpha\bar{\psi}^\alpha]. \quad (2.65)$$

The second term is often abbreviated as *c.c.* or *h.c.*, standing for the complex conjugate or the Hermitean conjugate.

2.7 The Lagrangian of the Standard Model

At this point, we can write down the Lagrangian describing a combined system gauge fields, scalar fields and fermions. We pick the gauge group $G = G_1 \times G_2 \times \dots$, and two-component fermions $\psi^{i\alpha}$ in some representation $\rho_{\text{fermion}}(g)_j^i$ of G , and scalars ϕ^a in some other representation $\rho_{\text{scalar}}(g)_b^a$ of G . Then the Lagrangian has the form

$$\begin{aligned} S = \int & \left[- \sum_s \left(\frac{1}{2g_{(s)}^2} \text{tr} F_{\mu\nu}^{(s)} F^{(s)\mu\nu} + \frac{\theta_{(s)}}{32\pi^2} \text{tr} F_{\mu\nu}^{(s)} F_{\rho\sigma}^{(s)} \epsilon^{\mu\nu\rho\sigma} \right) \right. \\ & - ((D_\mu\phi)^\dagger D^\mu\phi + V(\phi)) \\ & - (i\bar{\psi}\sigma^\mu D_\mu\psi + m_{ij}\psi^{i\alpha}\psi_\alpha^j + c.c.) \\ & \left. - (y_{aij}\phi^a\psi^{i\alpha}\psi_\alpha^i + c.c.) \right]. \end{aligned} \quad (2.66)$$

Here $g_{(s)}$ and $\theta_{(s)}$ are the coupling constants and the theta angles for the s -th gauge group G_s , $V(\phi)$ is a gauge-invariant scalar potential, m_{ij} is a gauge-invariant mass term, and y_{aij} specifies a gauge-invariant Yukawa interaction.

It so happened that when Vairocana Buddha³ decided to manifest itself as the real world, it chose the following G and the scalar/fermion representations.

³Please replace this entity by your favorite creator/universal entity.

2.7.1 Gauge group

The gauge group is

$$G = U(1) \times SU(2) \times SU(3). \quad (2.67)$$

(OK, it liked one, two, three. Why not.)

We denote its element by

$$g = (g_1, (g_2)_v^u, (g_3)_b^a) \quad (2.68)$$

where $u, v = 1, 2$ and $a, b = 1, 2, 3$ are indices for $SU(2)$ and $SU(3)$, respectively.

Here, $U(1)$ is the Hypercharge; $SU(2)$ is the Weak Force; and $SU(3)$ is the Strong Force. The electromagnetic $U(1)$ is the result of the Higgs mechanism acting on the hypercharge $U(1)$ and the $SU(2)$ weak force.

2.7.2 Scalar

Next, as the scalar, it chose $\phi^{u=1,2}$ with the transformation

$${}^g\phi^u = (g_1)^{+1/2} (g_2)_v^u \phi^v. \quad (2.69)$$

This is a charge $+1/2$ scalar in the fundamental two-dimensional representation of $SU(2)$, often denoted as **2** of $SU(2)$. (OK, that is one of the simplest nontrivial representation. Why not. But why only for $SU(2)$? And why does it have $U(1)$ hypercharge $+1/2$?) This is the Higgs field.

2.7.3 Fermions

Finally, as the fermions, it chose the following monstrosity.⁴ To describe them, let us first consider ψ_α consisting of six pieces,

$$\psi_\alpha = (Q_L, \quad \overline{u_R}, \quad \overline{d_R}, \quad \ell_L, \quad \overline{e_R}, \quad \overline{\nu_R})_\alpha, \quad (2.70)$$

where Q_L is the left-handed quark doublet, u_R is the right-handed up quark, d_R is the right-handed down quark, ℓ_L is the left-handed lepton doublet, e_R is the right-handed electron, ν_R is the right-handed neutrino. Here I used the complex conjugation to make everything left-handed.

They have the following gauge transformations:

$$\begin{aligned} {}^g(Q_L)^{ua} &= (g_1)^{+1/6} (g_2)_v^u (g_3)_b^a (Q_L)^{vb}, \\ {}^g(\overline{u_R})^{\bar{a}} &= (g_1)^{-2/3} \quad (\overline{g_3})_{\bar{b}}^{\bar{a}} (\overline{u_R})^{\bar{b}}, \\ {}^g(\overline{d_R})^{\bar{a}} &= (g_1)^{+1/3} \quad (\overline{g_3})_{\bar{b}}^{\bar{a}} (\overline{d_R})^{\bar{b}}, \\ {}^g(\ell_L)^u &= (g_1)^{-1/2} (g_2)_v^u \quad (\ell_L)^v, \\ {}^g\overline{e_R} &= (g_1)^{+1} \quad \overline{e_R}, \\ {}^g\overline{\nu_R} &= \quad \overline{\nu_R}. \end{aligned} \quad (2.71)$$

⁴Here I assume that the neutrino masses are given by the right-handed neutrinos. This part is not experimentally verified. Otherwise I need to add a non-renormalizable term in (2.66).

The same information is often tabulated as follows:

	U(1)	SU(2)	SU(3)
Q_L	+1/6	2	3
$\overline{u_R}$	-2/3	1	$\overline{\mathbf{3}}$
$\overline{d_R}$	+1/3	1	$\overline{\mathbf{3}}$
ℓ_L	-1/2	2	1
$\overline{e_R}$	+1	1	1
$\overline{\nu_R}$	0	1	1

(2.72)

or even

$$\begin{aligned} Q_L &: (\mathbf{2}, \mathbf{3})_{+1/6}, & \overline{u_R} &: (\mathbf{1}, \overline{\mathbf{3}})_{-2/3}, & \overline{d_R} &: (\mathbf{1}, \overline{\mathbf{3}})_{+1/3}, \\ \ell_L &: (\mathbf{2}, \mathbf{1})_{-1/2}, & \overline{e_R} &: (\mathbf{1}, \mathbf{1})_{+1}, & \overline{\nu_R} &: (\mathbf{1}, \mathbf{1})_0. \end{aligned} \quad (2.73)$$

And then introduce *three copies of it*:

$$\psi_\alpha^{i=1,2,3} = (Q_L, \overline{u_R}, \overline{d_R}, \ell_L, \overline{e_R}, \overline{\nu_R})_\alpha^{i=1,2,3}. \quad (2.74)$$

This additional index $i = 1, 2, 3$ labels the *flavor* and/or the *generations*.

Exercise 2.8. Isidore Rabi famously said “Who ordered that” when first member of the second generation, muon, was found. Locate the original reference to this quote.

2.7.4 Higgs potential, mass terms and Yukawa interactions

We need to specify the potential $V(\phi)$, the mass term $m\psi\psi$ and the Yukawa interaction $y\phi\psi\psi$ in (2.66), which should all be gauge-invariant. The potential up to the quartic term is fixed to have the form

$$V(\phi) = (\lambda/2)(\phi^\dagger\phi)^2 - m^2(\phi^\dagger\phi). \quad (2.75)$$

The quadratic term is known to be negative. We can only introduce gauge-invariant mass term to ν_R :

$$m\psi\psi + c.c. = m_{ij}^{\text{Maj.}} (\overline{\nu_R})^{i\alpha} (\overline{\nu_R})_\alpha^j + \overline{m}_{ij}^{\text{Maj.}} (\nu_R)_{\dot{\alpha}}^i (\nu_R)^{\dot{\alpha}j}. \quad (2.76)$$

The possible Yukawa interactions are of the form

$$y\phi\psi\psi = Y_{ij}^{\text{up}} \epsilon_{uv} \delta_{a\bar{a}} \phi^u (Q_L)^{iva\alpha} \overline{u_{R\alpha}}^{j\bar{a}} + Y_{ij}^{\text{down}} \delta_{\bar{u}v} \delta_{a\bar{a}} \overline{\phi}^{\bar{u}} (Q_L)^{iva\alpha} \overline{d_{R\alpha}}^{j\bar{a}} \quad (2.77)$$

$$+ Y_{ij}^{\text{lepton}} \epsilon_{uv} \phi^u (\ell_L)^{iv\alpha} \overline{e_{R\alpha}}^j + Y_{ij}^{\text{neutrino}} \delta_{\bar{u}v} \overline{\phi}^{\bar{u}} (\ell_L)^{iv\alpha} \overline{\nu_{R\alpha}}^j \quad (2.78)$$

Therefore, to completely specify the theory, we need to give the numerical values to three couplings and three theta angles

$$g_{1,2,3}, \quad \theta_{1,2,3}, \quad (2.79)$$

the two real parameters of the Higgs potential

$$\lambda, \quad m^2, \quad (2.80)$$

the Majorana neutrino masses

$$m_{ij}^{\text{Maj.}} \tag{2.81}$$

and the Yukawa couplings

$$Y_{ij}^{\text{up}}, \quad Y_{ij}^{\text{down}}, \quad Y_{ij}^{\text{lepton}}, \quad Y_{ij}^{\text{neutrino}}. \tag{2.82}$$

As we will see, some of parameters can be absorbed by field redefinitions.

2.8 Questions

We are then led to many natural questions:

- Why did the Buddha choose the parameters (2.79), (2.80), (2.81), (5.36) have the measured values?
- Why did the Buddha introduce three generations $i = 1, 2, 3$, not just one?
- Why did the Buddha decide that a single generation should consist of the strange representation (2.73)?
- Why did the Buddha choose the gauge group $U(1) \times SU(2) \times SU(3)$?
- Why did the Buddha use the Quantum Field Theory to implement us?

Exercise 2.9. Speculate why the Buddha made these decisions.

Exercise 2.10. If You are the One to create the universe, what would You do instead?

2.9 Anomaly cancellation

In fact, even the Buddha cannot choose an arbitrary fermion representation ρ_{fermion} , because of the following constraint. Let us say that the infinitesimal gauge transformation $g = e^{i\epsilon X}$ acts on the fermions ψ_{α}^i by

$${}^g\psi^i = \rho(g)_j^i \psi^j = \psi^i + i\epsilon \rho(X)_j^i \psi^j. \tag{2.83}$$

Then we have

Anomaly Cancellation Condition.

$\text{tr } \rho(X)^3$ and $\text{tr } \rho(X)$ need to vanish for arbitrary infinitesimal generator X .

We will see why this condition is necessary in later lectures. Very roughly, if this condition is not satisfied, the gauge invariance is broken, and one cannot contain the negative norm states in the system, and the positivity of the probability breaks down.

For simplicity, consider the case $G = U(1)$, with three fermions $\psi_\alpha^{i=1,2,3}$, with the charge q_1, q_2, q_3 respectively. The gauge transformation is

$$g\psi^1 = g^{q_1}\psi^1, \quad g\psi^2 = g^{q_2}\psi^2, \quad g\psi^3 = g^{q_3}\psi^3. \quad (2.84)$$

Equivalently, $\rho(g) = \text{diag}(g^{q_1}, g^{q_2}, g^{q_3})$. Writing $g = e^{i\epsilon X}$, we find

$$\rho(X) = \text{diag}(q_1, q_2, q_3)X. \quad (2.85)$$

Therefore the anomaly cancellation condition is

$$\text{tr } \rho(X)^3 = ((q_1)^3 + (q_2)^3 + (q_3)^3)X^3 = 0, \quad \text{tr } \rho(X) = (q_1 + q_2 + q_3)X = 0. \quad (2.86)$$

The general solution is

$$q_1 = -q_2, \quad q_3 = 0 \quad (2.87)$$

or its cyclic permutations.

Exercise 2.11. Show this.

Let us check that the anomaly cancellation condition is satisfied within a single generation (2.73). The gauge transformation parameter is $g = (g_1, g_2, g_3) \in U(1) \times SU(2) \times SU(3)$. The infinitesimal version can be parameterized as $g = e^{i\mathbb{X}}$ where

$$g_1 = e^{iX}, \quad g_2 = \text{diag}(e^{iY}, e^{-iY}), \quad g_3 = \text{diag}(e^{iA}, e^{iB}, e^{iC}) \quad (2.88)$$

where $A + B + C = 0$. $\rho(\mathbb{X})$ can be read off from the explicit gauge transformation rule (2.71) and we can compute $\text{tr } \rho(\mathbb{X})^3$ and $\text{tr } \rho(\mathbb{X})$. They are a cubic polynomial and a linear polynomial in X, Y , and A, B, C , respectively. And the anomaly cancellation condition demands that they are both zero.

Let us check this in the simpler case where $X \neq 0$ but $Y = A = B = C = 0$. We find

$$\text{tr } \rho(X)^3 = X^3 \left[\underbrace{3 \cdot 2 \cdot \left(\frac{1}{6}\right)^3}_{Q_L} + \underbrace{3 \cdot \left(-\frac{2}{3}\right)^3}_{\bar{u}_R} + \underbrace{3 \cdot \left(\frac{1}{3}\right)^3}_{\bar{d}_R} + \underbrace{2 \cdot \left(-\frac{1}{2}\right)^3}_{\ell_L} + \underbrace{1 \cdot 1^3}_{\bar{e}_R} \right] = 0. \quad (2.89)$$

Similarly, one finds

$$\text{tr } \rho(X) = X \left[\underbrace{3 \cdot 2 \cdot \left(\frac{1}{6}\right)}_{Q_L} + \underbrace{3 \cdot \left(-\frac{2}{3}\right)}_{\bar{u}_R} + \underbrace{3 \cdot \left(\frac{1}{3}\right)}_{\bar{d}_R} + \underbrace{2 \cdot \left(-\frac{1}{2}\right)}_{\ell_L} + \underbrace{1 \cdot 1}_{\bar{e}_R} \right] = 0. \quad (2.90)$$

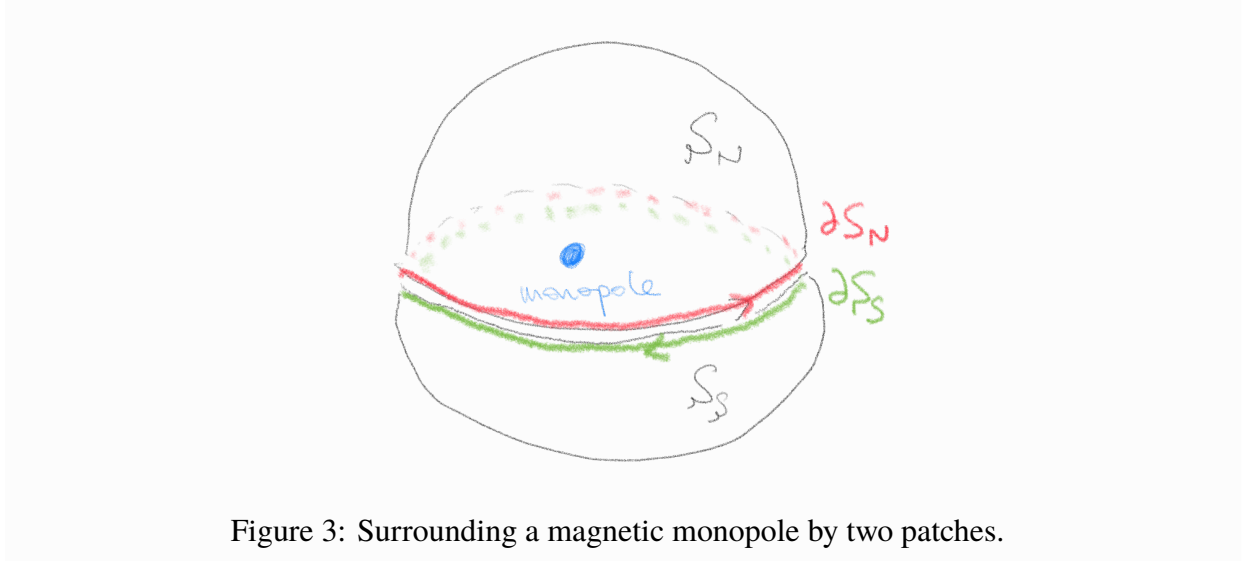


Figure 3: Surrounding a magnetic monopole by two patches.

Exercise 2.12. Check the cancellation of the anomaly in the most general case, without assuming $Y = A = B = C = 0$.

This motivates us to add the following question to those listed in Sec. 2.8:

- Why did the Buddha decide to satisfy the anomaly cancellation in this very strange manner?

3 Anomalies

3.1 Aharonov-Bohm phase and the Dirac quantization of monopole charges

Consider a 2d region S with a boundary ∂S . The Stokes theorem applied to $F_{\mu\nu} = \partial_\mu A_\nu - \partial_\nu A_\mu$ says that

$$\int_S F_{\mu\nu} dx^\mu dx^\nu = \int_{\partial S} A_\mu dx^\mu. \quad (3.1)$$

Note that the right hand side is gauge invariant, since the left hand side is. Note also that the right hand side can be nonzero even when $F_{\mu\nu}$ is zero at the boundary. The quantum-mechanical electron famously feels the phase

$$\exp(i \int_{\partial S} A_\mu dx^\mu), \quad (3.2)$$

known as the *Aharonov-Bohm phase* named after its discoverer; its effect was experimentally confirmed. In the math literature the same object is known as the *holonomy* of the gauge field.

Consider surrounding a magnetic monopole by two surfaces S_N and S_S , such that $C = \partial S_N = -\partial S_S$, see Fig. 3. The total magnetic flux is

$$\int_{S_N+S_S} F_{\mu\nu} dx^\mu dx^\nu = \int_C (A_N)_\mu dx^\mu - (A_S)_\mu dx^\mu = \int_C [(A_N)_\theta - (A_S)_\theta] d\theta. \quad (3.3)$$

The gauge fields A_N and A_S at the boundary is related by the gauge transformation

$$(A_N)_\theta = (A_S)_\theta + ig(\theta)\partial_\theta g(\theta)^{-1} \quad (3.4)$$

where $g(\theta)$ is a map from the circle parameterized by θ to the unit circle in the complex plane. Writing $g(\theta) = e^{i\varphi(\theta)}$ we find

$$\int_{S_N+S_S} F_{\mu\nu} dx^\mu dx^\nu = \int_C [(A_N)_\theta - (A_S)_\theta] d\theta = \int_C (\partial_\theta \varphi(\theta)) d\theta = 2\pi n \quad (3.5)$$

where n is the integer specifying the number of times the map $g(\theta)$ wraps around the target unit circle. This is the Dirac quantization of the magnetic monopole charge. Note that the same computation shows that

$$\exp(i \int_C (A_N)_\mu dx^\mu) = \exp(i \int_C (A_S)_\mu dx^\mu), \quad (3.6)$$

i.e. the phase felt by the electron around the monopole is a gauge-invariant quantity.

We will discuss 1+1d fermions below. The simplest nontrivial spacetime is the torus T^2 , given by identifying $t \sim t + T$ and $x \sim x + L$. The preceding argument can be carried out similarly, and we can show

$$\int dx dt F_{xt} = 2\pi n \quad (3.7)$$

for some integer $n \in \mathbb{Z}$. We will be mostly interested in the $n = 0$ case, with the further condition that $F_{\mu\nu} = 0$. This does not mean that the gauge field is trivial; there can be Aharonov-Bohm phases around the temporal and the spatial circle. For simplicity we take A_t and A_x are constants; this choice indeed leads to $F_{tx} = 0$. Still, we can have nonzero

$$g_t := \exp(i \int dt A_t) = e^{iT A_t}, \quad g_x := \exp(i \int dx A_x) = e^{iL A_x} \quad (3.8)$$

which are gauge-invariant. We note that the gauge transformation of the form

$$A_x \rightarrow A_x + g(x)\partial_x g(x)^{-1} \quad (3.9)$$

can change

$$\alpha := \int dx A_x \rightarrow \left(\int dx A_x \right) + 2\pi n' \quad (3.10)$$

for an integer n' , but the (exponentiated) Aharonov-Bohm phase g_x above is invariant. This transformation (3.10) is often called the large gauge transformation.

Let us see the effect of this spatial Aharonov-Bohm phase in the context of 1-dimensional quantum mechanics. The Hamiltonian of a free uncharged particle moving on a line with periodic boundary condition $x \sim x + L$ is $\frac{(i\partial_x)^2}{2m}$ acting on the wavefunction $\phi(x)$, which again has periodic boundary condition. Diagonalizing the Hamiltonian is easy; the eigenmodes and the energy eigenvalues are

$$\phi_n(x) = \exp\left(\frac{2\pi i n}{L} x\right); \quad E_n = \frac{1}{2m} \left(\frac{2\pi n}{L}\right)^2. \quad (3.11)$$

For a particle with charge q , the Hamiltonian is modified to

$$H = \frac{(iD_x)^2}{2m} = \frac{(i\partial_x + iqA_x)^2}{2m}. \quad (3.12)$$

The diagonalization is still easy:

$$\phi_n(x) = \exp\left(\frac{2\pi in}{L}x\right); \quad E_n = \frac{1}{2m} \left(\frac{2\pi n}{L} + A_x\right)^2 = \frac{1}{2m} \left(\frac{2\pi n + \alpha}{L}\right)^2. \quad (3.13)$$

We note that the large gauge transformation (3.10) does change the individual energy eigenvalues E_n , but the set $\{E_n\}$ is invariant under (3.10).

3.2 2d fermions

We will start our discussion of anomalies by first considering charged fermions in 1+1 dimensions. It is simpler than 3+1d fermions, and also has direct applications to the study of edge modes in the quantum Hall systems.

3.2.1 Massive fermion

A typical Lagrangian for 2d massive fermion has the form

$$S = \int dt dx (i\bar{\psi}_\ell(\partial_t - \partial_x)\psi_\ell + i\bar{\psi}_r(\partial_t + \partial_x)\psi_r + m\bar{\psi}_\ell\psi_r + \bar{m}\bar{\psi}_r\psi_\ell) \quad (3.14)$$

Indeed, by taking the variation, one finds

$$i(\partial_t - \partial_x)\psi_\ell = m\psi_r, \quad i(\partial_t + \partial_x)\psi_r = \bar{m}\psi_\ell. \quad (3.15)$$

Expanding in Fourier modes $\exp(iEt - ipx)$, we find

$$-(E + p)\psi_\ell = m\psi_r, \quad -(E - p)\psi_r = \bar{m}\psi_\ell \quad (3.16)$$

from which we find $E^2 - p^2 = |m|^2$.

Exercise 3.1. Derive the 2d fermion Lagrangian in the usual manner, by constructing the γ matrices. Rewrite the resulting Lagrangian into the form (3.14).

3.2.2 Massless fermion and its gravitational anomaly

A special case is when $m = 0$, where we have $E^2 = p^2$, which can be solved as $E = \pm p$ in 1+1d. Correspondingly, we can consider a left-moving or a right-moving fermion separately. For example, a purely left-moving fermion has the Lagrangian

$$S = \int dt dx i\bar{\psi}_\ell(\partial_t - \partial_x)\psi_\ell \quad (3.17)$$

whose classical equation of motion is

$$i(\partial_t - \partial_x)\psi_\ell = 0 \quad (3.18)$$

whose general solution is

$$\psi_\ell = f(x + t). \quad (3.19)$$

Let us analyze its quantization, under the periodic boundary condition $x \sim x + L$. We expand ψ_ℓ into Fourier modes:

$$\psi_\ell(x) = \sum_n \psi_n \exp\left(\frac{2\pi i n}{L}x\right). \quad (3.20)$$

Then the anticommutation relation is

$$\{\psi_n, \psi_m^\dagger\} = \delta_{nm}. \quad (3.21)$$

The mode ψ_n carries the momentum $-2\pi n/L$ and the energy $2\pi n/L$. Considered in isolation, this leads to two states

$$|\downarrow\rangle_n = \psi_n^\dagger |\uparrow\rangle_n, \quad |\uparrow\rangle_n = \psi_n |\downarrow\rangle_n. \quad (3.22)$$

It is then natural to assign the momentum $\mp \frac{1}{2}(2\pi n/L)$ and the energy $\pm \frac{1}{2}(2\pi n)/L$ to the states $|\uparrow\rangle_n$ and $|\downarrow\rangle_n$.

The vacuum state is then obtained by using $|\downarrow\rangle_n$ for $n > 0$ and $|\uparrow\rangle_n$ for $n < 0$. The energy and the momentum of the vacuum are both given by

$$E = -P = -\frac{2\pi}{L} \sum_{n \geq 1} n = \frac{1}{12} \frac{2\pi}{L}. \quad (3.23)$$

The Casimir energy has the sign opposite to that of a scalar (1.56). What is more strange is the Casimir momentum. Recall that e^{iPX} corresponds to the shift of the x direction by X . Since we started from a periodic boundary condition $x \sim x + L$, one might expect that the shift by $X = L$ should not change the system. But we gain the phase

$$\exp\left(i \frac{1}{12} \frac{2\pi}{L} L\right) = \exp\left(\frac{2\pi i}{12}\right). \quad (3.24)$$

This is weird, and is known as the gravitational anomaly.

An overall shift of the periodic direction does nothing classically, but the quantum system responds nontrivially. Such phenomena are called quantum anomalies. In this case the geometric shift is a baby version of the general coordinate transformation, which is why this particular case is known as the gravitational anomaly. A system with gravitational anomaly cannot couple consistently with quantum gravity: the negative norm states of the graviton cannot be consistently decoupled.

By a long chain of arguments which I do not have time to develop here, the gravitational anomaly of the 1+1d boundary system is known to be proportional to the thermal Hall conductivity of the 2+1 dimensional bulk. So this is an experimentally measurable quantity and has indeed been measured.

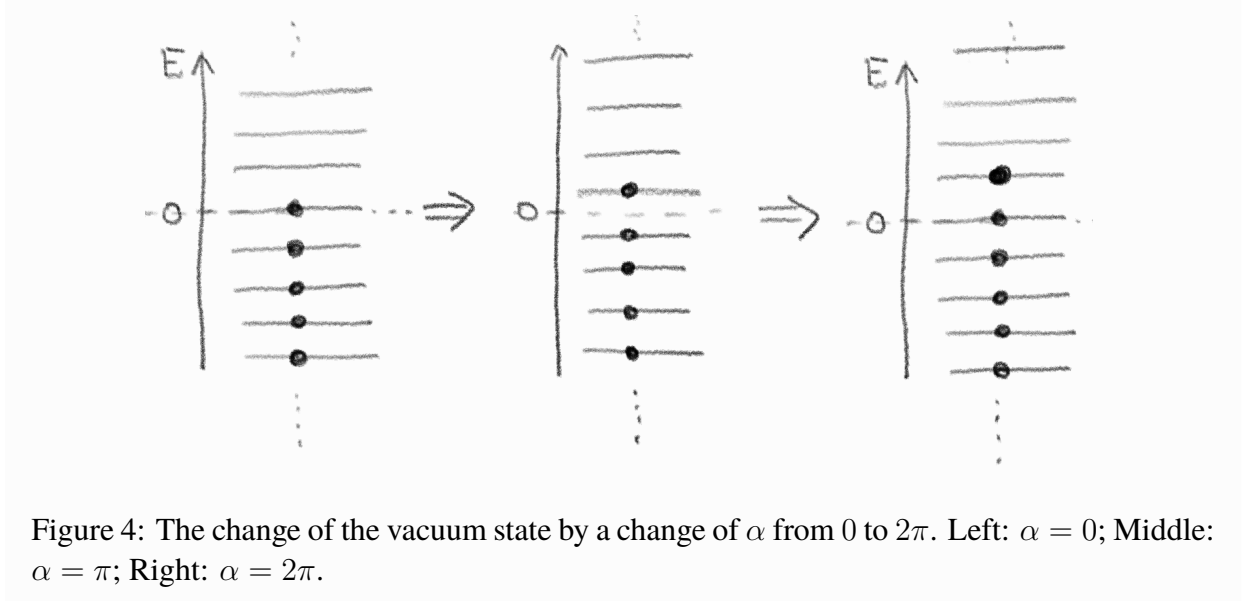


Figure 4: The change of the vacuum state by a change of α from 0 to 2π . Left: $\alpha = 0$; Middle: $\alpha = \pi$; Right: $\alpha = 2\pi$.

Note that the right-moving fermion has

$$E = +P = \frac{1}{12} \frac{2\pi}{L}. \quad (3.25)$$

Therefore, if we have both a right-moving and a left-moving fermion, the Casimir energies sum up, but the gravitational anomalies cancel out. In particular, a massive fermion does not have the gravitational anomaly.

3.2.3 Charged massless fermion and its gauge anomaly

Let us next consider the case of a charged massless fermion

$$S = \int dt dx i \bar{\psi}_\ell (D_t - D_x) \psi_\ell \quad (3.26)$$

with a periodic boundary condition $x \sim x + L$ with a spatial Aharonov-Bohm phase $\alpha = A_x L$. We first consider the case when the gauge transformation is ${}^g\psi_\ell = g\psi_\ell$, i.e. when ψ_ℓ is of charge 1. Each mode now has the momentum $-(2\pi n + \alpha)/L$ and the energy $(2\pi n + \alpha)/L$.

When $|\alpha|$ is very small, the vacuum state is still obtained by using $|\downarrow\rangle_n$ for $n > 0$ and $|\uparrow\rangle_n$ for $n < 0$.

Exercise 3.2. Compute the Casimir energy and the Casimir momentum when $\alpha \neq 0$, by the regularization as in Sec. 1.7.

However, when we gradually increase α from 0 to 2π , we have the situation shown in Fig. 4. As each mode carries charge +1, we see that there is a gain of charge +1 in this process. It is

natural to assign the charge of the vacuum to be

$$Q = \frac{\alpha}{2\pi}. \quad (3.27)$$

Exercise 3.3. Derive this vacuum charge by regularizing the charge of the vacuum state as in the case of the Casimir energy.

This is known as the gauge anomaly: two situations $\alpha = 0$ and $\alpha = 2\pi$ are gauge equivalent and should give rise to the same physical quantity, if the system is invariant under the gauge transformation. But here we clearly have two different answers for these two situations.

We can generalize the analysis to the case of charge q fermion. The covariant derivative is $i\partial_x + iqA_x$, and therefore the energy eigenvalues of the modes are $(2\pi n + q\alpha)/L$. Then the shift $\alpha \rightarrow \alpha + 2\pi$ gives us q additional filled states. Each mode carries charge q , and therefore we gained the charge q^2 in total. We then have

$$Q = \frac{q^2\alpha}{2\pi} \quad (3.28)$$

as the vacuum gauge charge.

We can also consider the right-moving fermions. It has the energy $-(2\pi n + q\alpha)/L$ and thus has the opposite shift of the charge, with the anomaly

$$Q = -\frac{q^2\alpha}{2\pi}. \quad (3.29)$$

In general, we can consider the system where we have left-moving fermions of charge q_1, q_2, \dots and right-moving fermions of charge $\tilde{q}_1, \tilde{q}_2, \dots$. The total anomaly is

$$Q = \frac{\alpha}{2\pi} [\sum (q_i)^2 - \sum (\tilde{q}_i)^2]. \quad (3.30)$$

Therefore, the gauge anomaly is absent only when

$$\sum (q_i)^2 - \sum (\tilde{q}_i)^2 = 0. \quad (3.31)$$

This is the anomaly cancellation condition in 1+1d.

We note that a massive fermion such as (3.14) contains a left-moving component of charge $+q$ and a right-moving component of the same charge $+q$. Therefore their contribution to the anomaly automatically cancels.

3.3 An index theorem

We switch gears and ask the following question [AC79]. Let us consider a non-relativistic charged particle moving in 2d under the effect of the magnetic field $B = F_{xy}$. The Hamiltonian is (up to a

proportionality factor)

$$H = (iD_x)^2 + (iD_y)^2 + B(x, y)\sigma_z = \quad (3.32)$$

$$= \begin{pmatrix} i(D_x - iD_y)i(D_x + iD_y) & 0 \\ 0 & i(D_x + iD_y)i(D_x - iD_y) \end{pmatrix} \quad (3.33)$$

$$= \begin{pmatrix} 0 & i(D_x + iD_y)^2 \\ i(D_x - iD_y) & 0 \end{pmatrix}. \quad (3.34)$$

We are interested in its spectrum.

Setting $D := i(D_x + iD_y)$, this falls within a more general problem of the diagonalization of

$$D^\dagger D|\psi_+\rangle = E_+|\psi_+\rangle, \quad DD^\dagger|\psi_-\rangle = E_-|\psi_-\rangle. \quad (3.35)$$

Let us first study this. We note that

$$\langle\psi_+|E_+|\psi_+\rangle = \langle\psi_+|D^\dagger D|\psi_+\rangle = \|D|\psi_+\rangle\|^2. \quad (3.36)$$

Therefore,

$$D|\psi_+\rangle = 0 \iff E_+ = 0 \quad (3.37)$$

and otherwise $E_+ > 0$. Similarly,

$$D^\dagger|\psi_-\rangle = 0 \iff E_- = 0 \quad (3.38)$$

and otherwise $E_- > 0$.

Now suppose we found an eigenstate $|\psi_+\rangle$ of $D^\dagger D$ with eigenvalue $E > 0$ such that $\langle\psi_+|\psi_+\rangle = 1$. Let $|\psi_-\rangle := D|\psi_+\rangle/\sqrt{E}$. This satisfies $\langle\psi_-|\psi_-\rangle = 1$ and

$$DD^\dagger|\psi_-\rangle = DD^\dagger D \frac{|\psi_+\rangle}{\sqrt{E}} = DE_+ \frac{|\psi_+\rangle}{\sqrt{E_+}} = E|\psi_-\rangle. \quad (3.39)$$

This means that $|\psi_-\rangle$ thus defined is an eigenstate of DD^\dagger with the same eigenvalue. We also have $|\psi_+\rangle = D^\dagger|\psi_-\rangle/\sqrt{E}$.

Therefore, eigenstates of zero energy of $D^\dagger D$ and DD^\dagger are annihilated by D and D^\dagger respectively, and eigenstates of nonzero energy of $D^\dagger D$ and DD^\dagger form pairs with the same eigenvalues. We denote the eigenstates of zero energy by $|\Psi_+^u\rangle$ and $|\Psi_-^s\rangle$, and the eigenstates of nonzero energy E_a by $|\psi_+^a\rangle$ and $|\psi_-^a\rangle$, where $|\psi_-^a\rangle = D|\psi_+^a\rangle$ and $|\psi_+^a\rangle = D^\dagger|\psi_-^a\rangle$.

Let us come back to our case when $D = i(D_x - iD_y)$, and study the eigenstates of zero energy more carefully. They are often called zero modes.

Recall that $B = \partial_x A_y - \partial_y A_x$. We further write $A_x = -\partial_y \rho$ and $A_y = +\partial_x \rho$. Then $B = (\partial_x^2 + \partial_y^2)\rho$, which can be solved. We now define $f_+(x, y)$ by

$$\Psi_+(x, y) = f_+(x, y)e^{\rho(x, y)}. \quad (3.40)$$

Then

$$(D_x - iD_y)\Psi_+(x, y) = 0 \iff (\partial_x - i\partial_y)f_+(x, y) = 0. \quad (3.41)$$

which means that $f_+(x, y)$ is a holomorphic function of $w = x + iy$. Similarly, by defining

$$\Psi_-(x, y) = f_-(x, y)e^{-\rho(x, y)}. \quad (3.42)$$

one finds

$$(D_x + iD_y)\Psi_-(x, y) = 0 \iff (\partial_x + i\partial_y)f_-(x, y) = 0. \quad (3.43)$$

which means that $f_-(x, y)$ is an antiholomorphic function, i.e. it depends on \bar{w} but not on w .

Consider in particular that B is nonzero only in a finite region $S \subset \mathbb{R}^2$ and zero outside. Asymptotically, we have

$$\rho(x, y) \sim \frac{1}{2\pi} \left(\int_S B dx dy \right) \log r = n \log r \quad (3.44)$$

where we write $\int_S B dx dy = 2\pi n$ and r is the radial distance from a fixed point in the region S . For definiteness let $n > 0$, and for simplicity we assume n is an integer. This means that

$$e^{+\rho} \sim r^n, \quad e^{-\rho} \sim r^{-n} \quad (3.45)$$

asymptotically. We require that $|\Psi_\pm\rangle$ to decay when $r \rightarrow 0$, so that it is square integrable. This means that $f_+(x, y)$ is forced to be zero, while $f_-(x, y)$ is a linear combination of \bar{w}^0 to \bar{w}^{n-1} .

We conclude that there are no zero energy states $|\Psi_+\rangle$, and there are n zero energy states $|\Psi_-^s\rangle$, $s = 1, \dots, n$. When we have a charge q particle, we basically have the same conclusion, except that $s = 1, \dots, qn$.

3.4 Anomalies of 4d chiral fermion

Let us now consider a charge $q > 0$ chiral fermion in 3+1d. The Lagrangian is

$$S = i \int d^4x i \bar{\psi}^{\dot{\alpha}} \sigma_{\dot{\alpha}\alpha}^\mu D_\mu \psi^\alpha \quad (3.46)$$

$$= i \int dt dz dx dy (\bar{\psi}^1, \bar{\psi}^2) \begin{pmatrix} D_t + D_z & D_x + iD_y \\ D_x - iD_y & D_t - D_z \end{pmatrix} (\psi^1, \psi^2). \quad (3.47)$$

We now introduce $B_z(x, y) = F_{xy}$ in a finite region $S \subset \mathbb{R}_{x,y}^2$ and zero outside. We further put a periodic boundary condition $z \sim z + L$. We first expand the chiral fermions using the eigenmodes of H studied in the last section:

$$\psi^1(z, t; x, y) = \sum_a \psi_r^a(z, t) |\psi_+^a\rangle, \quad (3.48)$$

$$\psi^2(z, t; x, y) = \sum_a \psi_\ell^a(z, t) |\psi_-^a\rangle + \sum_{s=1}^{qn} \Psi_\ell^s |\Psi_-^s\rangle. \quad (3.49)$$

Plugging this into (3.47), we obtain

$$S = \sum_{s=1}^{qn} \int dt dx (i\bar{\psi}_\ell^s (D_t - D_x) \psi_\ell^s) + \sum_a \int dt dx (i\bar{\psi}_\ell^a (D_t - D_x) \psi_\ell^a + i\bar{\psi}_r^a (D_t + D_x) \psi_r^a + \sqrt{E_a} \bar{\psi}_\ell^a \psi_r^a + \sqrt{E_a} \bar{\psi}_r^a \psi_\ell^a). \quad (3.50)$$

This counts as qn left-moving 2d fermions of charge q , thus contributing

$$Q = \frac{q^3 n \alpha}{2\pi} \quad (3.51)$$

to the anomalous gauge charge, and

$$P = \frac{qn}{12} \frac{2\pi}{L} \quad (3.52)$$

to the gravitational anomaly. When $q < 0$, we instead have $|q|n$ right-moving 2d fermions, contributing to

$$Q = -\frac{|q|^3 n \alpha}{2\pi} = \frac{q^3 n \alpha}{2\pi} \quad (3.53)$$

to the anomalous gauge charge, and

$$P = -\frac{q|n|}{12} \frac{2\pi}{L} = \frac{qn}{12} \frac{2\pi}{L} \quad (3.54)$$

to the gravitational anomaly.

When there are multiple chiral 4d fermions of charge q_a , we therefore see that the total gauge anomaly is proportional to

$$\sum_a (q_a)^3, \quad (3.55)$$

whereas the total gravitational anomaly is proportional to

$$\sum_a q_a. \quad (3.56)$$

When we treat $U(1)$ gauge symmetry quantum mechanically, the gauge anomaly has to vanish, in order to contain the negative norm state of the photons. When we treat the general coordinate transformation invariance quantum mechanically, the gravitational anomaly has to vanish, in order to contain the negative norm state of the gravitons.

Exercise 3.4. Here I chose a rather non-standard method to derive the anomaly. Learn more standard methods in your favorite textbook.

4 Spontaneous symmetry breaking, Higgs effect and solitons

The topic in this section is the Higgs effect and the associated solitons. Everything stated in this section concerns classical field theory, and no quantum physics is involved.



Figure 5: The form of the potentials $V(\phi)$ $m^2 > 0$ and $m^2 < 0$.

4.1 Spontaneous symmetry breaking

Let us start by considering a complex scalar field ϕ with the following action:

$$S = \int d^4x [-\partial^\mu \bar{\phi} \partial_\mu \phi - V(\phi)], \quad V(\phi) = m^2 |\phi|^2 + \frac{\lambda}{2} |\phi|^4. \quad (4.1)$$

The action is invariant under the $U(1)$ transformation

$${}^g \phi(x) = g \phi(x), \quad (4.2)$$

where $g \in U(1)$ is a complex number of absolute value 1, which is *independent* of the spacetime.

We require $\lambda > 0$. Otherwise, the potential is unbounded below when $|\phi|$ is sufficiently large, which is very dangerous. We do not require m^2 to be positive. See Fig. 5 for a crude drawing of the shape of the potential. The case $m^2 < 0$ is sometimes known as the Mexican hat potential or the wine-bottle potential. The latter nomenclature is somewhat common in the Japanese particle physics literature but is rarely found outside.

Exercise 4.1. Track down who first used these two suggestive names for the shape of the potential.

Let us consider the spacetime-independent situation, specified by giving a value to the scalar field ϕ . This is usually denoted by $\langle \phi \rangle$, and called the vacuum expectation value, or *vev* for short. The equation of motion says that $\partial V / \partial \phi = 0$ when evaluated at $\phi = \langle \phi \rangle$.

When $m^2 > 0$, there is a single solution $\langle \phi \rangle = 0$. When $m^2 < 0$, the choice $\langle \phi \rangle = 0$ is still a solution, but

$$\langle \phi \rangle = \frac{|m|}{\sqrt{\lambda}} e^{i(\theta)} \quad (4.3)$$

for an arbitrary $\langle \theta \rangle$ also forms a continuous family of solutions. The notation

$$v := |\langle \phi \rangle| = \frac{|m|}{\sqrt{\lambda}} \quad (4.4)$$

is often also found.

When $\langle \phi \rangle = 0$, the vacuum expectation value is invariant under the phase rotation (4.2). This vacuum is said to preserve the U(1) symmetry. Equivalently, the U(1) symmetry is unbroken there. Let us analyze the small variations $\delta\phi(t, x, y, z)$ around it. Considering the mode $\propto e^{i(Et - \vec{p}\vec{x})}$, we see that we have

$$E^2 - p^2 = m^2. \quad (4.5)$$

When $m^2 > 0$, this simply says that (after quantization) the particle has the mass $|m|$. The $\lambda|\phi|^4/2$ term then provides an interaction term. When $m^2 < 0$, we see that the 4-momentum is spacelike; naively, the speed of the excitation exceeds the speed of light. Such a mode is called tachyonic. A better interpretation is obtained by setting $p = 0$, $iEt = +|m|t$ instead, getting the behavior $\delta\phi \sim e^{+|m|t}$. This means that the value of the scalar field starts to exponentially grow. This behavior is clear from the potential shown in Fig. 5: to save energy, the field rolls down the potential to the bottom (4.3).

There, the U(1) symmetry action (4.2) does change the value of $\langle \phi \rangle$. In this situation, the U(1) symmetry is said to be broken. Writing $\phi = \rho e^{i\theta}$, we can consider the fluctuation $\delta\rho$ and $\delta\theta$ around it. We see that $\delta\rho$ has the mass

$$m_\rho^2 = \frac{1}{2} \frac{\partial^2}{\partial\rho\partial\rho} V \Big|_{\phi=e^{i\theta}\sqrt{\lambda}/|m|} = 2|m^2|, \quad (4.6)$$

whereas $\delta\theta$ is massless,

$$m_\theta^2 = 0. \quad (4.7)$$

Exercise 4.2. Confirm these two masses.

It is a general feature of a relativistic field theory that for each broken continuous symmetry, there is a massless particle. This is the theorem of Nambu and Goldstone; the corresponding particle is called a Nambu-Goldstone mode. At this classical level, this is simply because $\langle \phi \rangle$ and $g\langle \phi \rangle$ has the same energy $V(\phi)$. So there is no energy variation in the direction obtained by an infinitesimal symmetry operation applied to $\langle \phi \rangle$, meaning that this direction is massless.

We can take a limit where we keep $v = |m|/\lambda$ fixed but sending both $|m|$ and λ infinity. Then the radial mode ρ is infinitely massive, and is decoupled from the rest. The Lagrangian for the remaining mode θ is simply

$$S = \int d^4x |v|^2 (-\partial^\mu \theta \partial_\mu \theta), \quad (4.8)$$

where we need to remember that $\theta \sim \theta + 2\pi$.

4.2 Higgs mechanism in U(1) gauge theory

We now couple the Maxwell field to the model studied above:

$$S = - \int d^4x \left[\frac{1}{4e^2} F_{\mu\nu} F^{\mu\nu} + \overline{D^\mu \phi} D_\mu \phi + V(\phi) \right], \quad V(\phi) = m^2 |\phi|^2 + \frac{\lambda}{2} |\phi|^4 \quad (4.9)$$

where

$$D_\mu \phi = (\partial_\mu + iA_\mu)\phi. \quad (4.10)$$

This model is now invariant under the spacetime-dependent transformation

$${}^g A_\mu = A_\mu - ig\partial_\mu g^{-1}, \quad {}^g \phi = g\phi. \quad (4.11)$$

When $m^2 > 0$, the only vacuum (i.e. the configuration with the lowest energy) is $\langle \phi \rangle = 0$. The field ϕ and A_μ are both massless. This is called the Coulomb phase.

When $m^2 < 0$, the extrema of $V(\phi)$ are $\langle \phi \rangle = 0$ and

$$\langle \phi \rangle = |m|/\sqrt{\lambda}e^{i\theta} \quad (4.12)$$

as before. The first choice $\langle \phi \rangle = 0$ is unstable, since the field ϕ is tachyonic. We concentrate on the latter.

We recall that the gauge transformation $g = e^{i\chi}$ shifts the phase of ϕ :

$${}^g \langle \theta \rangle = \langle \theta \rangle + \chi. \quad (4.13)$$

In the non-gauge version we studied in the last subsection, this meant that there are continuous family of vacua related by the U(1) global symmetry. In the gauged version here, the two configurations related by the gauge symmetry is *identified* instead. Therefore, the family (4.12) is in fact a single configuration.

Writing $\phi = \rho e^{i\theta}$ as before, we see that the mode $\delta\rho$ still has the mass $m_\rho^2 = |m|^2/\lambda$. The mode $\delta\theta$ is a gauge mode and is gone. Let us now analyze the gauge field. We note that $D_\mu \phi = (\partial_\mu + iA_\mu)\phi$ now contains the piece $iA_\mu \langle \phi \rangle$. This means that we have

$$\frac{1}{4e^2} F_{\mu\nu} F^{\mu\nu} + \overline{D^\mu \phi} D_\mu \phi \supset \frac{1}{4e^2} F_{\mu\nu} F^{\mu\nu} + A_\mu A^\mu |v|^2. \quad (4.14)$$

The equation of motion is then

$$\partial^\mu F_{\mu\nu} = 2e^2 v^2 A_\nu. \quad (4.15)$$

Applying ∂^ν on both sides, we see $\partial^\nu A_\nu = 0$. Plugging this back to the left hand side and replacing ∂^μ by ik^μ , we find that $k^2 = 2e^2 v^2$, i.e. the mass of the gauge boson is now

$$m_A^2 = 2e^2 v^2 = 2e^2 |m^2|/\lambda. \quad (4.16)$$

We recall that the mass of $\delta\rho$ was found above to be

$$m_\rho^2 = 2|m^2|. \quad (4.17)$$

This is the Higgs mechanism in the U(1) gauge theory. In this context, the field ϕ is called the Higgs field, and the radial mode $\delta\rho$ is called the Higgs boson. This phase is known as the Higgs phase.

We note that a massive gauge boson has three polarizations. Taking $k^\mu = (E, 0, 0, 0)$, we find $A_t = 0$ and $A_{x,y,z}$ can be nonzero. Recall that a massless gauge boson had two polarizations:

taking $k^\mu = (E, E, 0, 0)$, $A_t + A_x = 0$, $(A_t, A_x) \propto (1, -1)$ were gauge modes and nonphysical, and the two polarizations came from A_y and A_z . The difference in the number of polarizations is often said to come from the gauge field A_μ eating the Nambu-Goldstone mode $\delta\theta$.

We can again take a limit where we keep $v = |m|/\lambda$ fixed but sending both $|m|$ and λ infinity. We then force $|\phi| = v$ and we can write $\phi = ve^{i\theta}$ so that

$$\overline{D_\mu \phi} D^\mu \phi = |v|^2 (\partial_\mu \theta + A_\mu) (\partial^\mu \theta + A^\mu). \quad (4.18)$$

The Lagrangian of the system is now

$$S = - \int d^4x \left[\frac{1}{4e^2} F_{\mu\nu} F^{\mu\nu} + |v|^2 (\partial_\mu \theta + A_\mu) (\partial^\mu \theta + A^\mu) \right], \quad (4.19)$$

and describes a single massive vector boson. This simplified form of the Higgs mechanism is called the Stückelberg mechanism, and θ the Stückelberg field.

4.3 Under the external magnetic field

Let us consider applying an external magnetic field $B_{\text{ext}} = F_{xy}$ to the system in the Higgs phase. We consider time-independent configurations. The energy to be minimized is

$$\int d^3x \left[\frac{1}{2e^2} (B - B_{\text{ext}})^2 + |D_i \phi|^2 + V(\phi) \right]. \quad (4.20)$$

We note that this has basically the same form as the energy functional of the Ginzburg-Landau effective model of the superconductor, except that the field ϕ has charge $q = 2$ instead of charge $q = 1$ as assumed here. This is due to the fact that in a superconducting material, what gets the vev is the Cooper pair, which is a combination of two electrons. In any case, the physics we see below is essentially the same, and we borrow the terminology from there.

Let us first consider spatially independent situation. One is to take $\langle \phi \rangle = 0$ and $B_{\text{ext}} = B$. The energy density is zero. Another is to take $D_i \phi = 0$ and $\langle \phi \rangle = |v|$. We note that $F_{ij} \phi = i[D_i, D_j] \phi = 0$, meaning that $B = 0$. Therefore the energy density is

$$\frac{1}{2e^2} B_{\text{ext}}^2 + V(v) = \frac{1}{2e^2} B_{\text{ext}}^2 - \frac{m^4}{2\lambda}. \quad (4.21)$$

This means that when

$$B_{\text{ext}} < \frac{e|m^2|}{\sqrt{\lambda}} =: B_c, \quad (4.22)$$

the external magnetic field is screened. This is called the Meissner effect. In contrast, when $B_{\text{ext}} > B_c$, the Coulomb phase with $\langle \phi \rangle = 0$ is favored.

So far we only considered spatially independent configurations. When $B_{\text{ext}} = 0$, the origin $\langle \phi \rangle$ contains a tachyon and is unstable. Let us see what happens with $B_{\text{ext}} \neq 0$. The equation to be considered is

$$-\partial_i^2 \phi = [-(D_x^2 + D_y^2 + D_z^2) - |m^2|] \phi = 0. \quad (4.23)$$

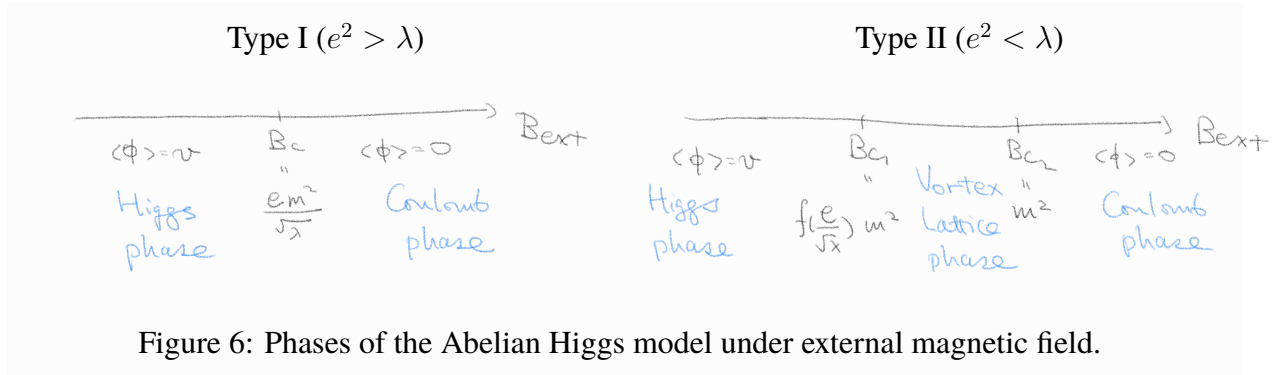


Figure 6: Phases of the Abelian Higgs model under external magnetic field.

We recall that $[D_x, D_y] = iB_{\text{ext}}$, and therefore the smallest eigenvalue of $D_x^2 + D_y^2$ is B_{ext} . This means that when

$$B_{\text{ext}} > m^2, \quad (4.24)$$

the Coulomb phase is stable. Therefore, when $e^2 > \lambda$, $B_c > m^2$ and the Coulomb phase is indeed stable. However, when $e^2 < \lambda$, the spatially independent solution is unstable. What will happen in this case?

We will see below that we will have a lattice of vortices. Before giving the detail, a summary of the phases are given in Fig. 6. Following the terminology in the study of superconductors, the case $e^2 > \lambda$ and the case $e^2 < \lambda$ are known as Type I and Type II, respectively.

4.4 Vortex solution

Let us next consider $|v|$ is space-dependent and is small in a finite region C in the xy plane around $x = y = 0$. One can then turn on the magnetic field in C without costing too much energy. What is the situation outside C ?

Consider a large circle of radius r around the region C , parameterized by $0 < \varphi < 2\pi$. Recall the integral $\int d\varphi \partial_\varphi \theta = 2\pi n$ for an integer n , which counts the winding number of the field $\theta(\varphi)$ as we traverse the circle by varying φ . By a gauge transformation we can arrange A_φ to be constant in φ . Now we would like to minimize the energy:

$$E \geq \int_{r>R} r dr \int d\varphi |v|^2 \frac{1}{r} (\partial_\varphi \theta + A_\varphi)^2 \quad (4.25)$$

The minimum is obtained when $\partial_\varphi \theta$ is also constant, i.e. when $\theta = n\varphi$. We then have

$$E \geq \int_{r>R} \frac{dr}{r} |v|^2 |n + A_\varphi|^2. \quad (4.26)$$

This diverges logarithmically unless A_φ is an integer $-n$. Recalling that $\int_{\partial C} A_\varphi d\varphi = \int_C F_{xy} dx dy$, we see that

$$\frac{1}{2\pi} \int_C F_{xy} dx dy = -n \in \mathbb{Z}. \quad (4.27)$$

This shows that the magnetic flux in the Higgs phase is quantized, and is proportional to the winding number of the field θ . This is known as a Abrikosov-Nielsen-Olsen vortex.⁵

Exercise 4.3. Redo the analysis when the original scalar field $\phi = \rho e^{i\theta}$ has charge $q \in \mathbb{Z}$ instead of 1 as assumed above. Show that $\frac{1}{2\pi} \int_C F_{xy} dx dy$ is now of the form n/q for $n \in \mathbb{Z}$. In the normal superconductor, $q = 2$ since ϕ is the field representing the Cooper pair.

We can see that a vortex has a fixed size roughly as follows. Let us say we put a magnetic field $\int dx dy B = 2\pi n$ in the region C of area A . To do this, we set $\langle \phi \rangle = 0$ in the same region. The energy const is then

$$\int dx dy \left[\frac{1}{2e^2} B^2 + (V(0) - V(v)) \right] \sim \frac{(2\pi n)^2}{A} + (V(0) - V(v))A. \quad (4.28)$$

This has a minimum at a specific value A .

The case $|n| = 1$ is the minimal vortex. Whether the case $|n| = 2$ splits into two minimal vortices or prefers to be a single vortex depends on the parameter e^2/λ . The fact that Type II admits a vortex phase means that the vortices prefer to separate in $e^2 < \lambda$ whereas the vortices clump together when $e^2 > \lambda$.

Exercise 4.4. Come up with a more direct argument showing that the vortices repel each other when $e^2 \ll \lambda$ and clump together when $e^2 \gg \lambda$.

For the type II system, when the external magnetic field B_{ext} is above B_{c1} the magnetic field penetrates the system as a combination of minimal vortices, each having $\int B dx dy = 2\pi$. They are packed together to realize the given external magnetic field B_{ext} . As we increase B_{ext} , the vortices become more and more tightly packed, and at $B_{\text{ext}} = B_{c2}$ they all coalesce and become uniform. From the discussion in the last subsection, we know $B_{c2} = m^2$. From dimensional analysis and scaling argument, B_{c1} is of the form $B_{c1} = f(e/\sqrt{\lambda})m^2$ where $f(x)$ is a dimensionless function.

Exercise 4.5. Learn how to compute $f(x)$.

4.5 Bogomolny trick

When $e^2 = \lambda$, we can use the following trick first found by Bogomolny to study the system in detail. Let us assume that the system is independent of t and z . The energy to be minimized is

$$E := \int dx dy \left[\frac{1}{2e^2} B^2 + \overline{D_x \phi} D_x \phi + \overline{D_y \phi} D_y \phi + \frac{e^2}{2} (|\phi|^2 - v^2)^2 - \frac{m^4}{2\lambda} \right] \quad (4.29)$$

⁵The model without the gauge field, with logarithmically divergent energy, was considered by Nielsen and Olsen. Abrikosov found the finite-energy solution with the gauge field.

where we recall $v^2 = |m^2|/\lambda$. We first note that

$$\frac{1}{2e^2}B^2 + \frac{e^2}{2}(|\phi|^2 - \frac{|m^2|}{\lambda})^2 = \frac{1}{2e^2}(B + e^2(|\phi|^2 - v^2))^2 - B(|\phi|^2 - v^2). \quad (4.30)$$

We then consider

$$\overline{D_x\phi}D_x\phi + \overline{D_y\phi}D_y\phi = (D_x + iD_y)\overline{\phi}(D_x - iD_y)\phi + i(D_x\overline{\phi}D_y\phi - D_y\overline{\phi}D_x\phi) \quad (4.31)$$

We perform a partial integral of the second term and note

$$\int dx dy [\overline{D_x\phi}D_x\phi + \overline{D_y\phi}D_y\phi] = \int dx dy [(D_x - iD_y)\phi]^2 - i(\overline{\phi}(D_xD_y - D_yD_x)\phi) \quad (4.32)$$

$$= \int dx dy [(D_x - iD_y)\phi]^2 + B|\phi|^2. \quad (4.33)$$

Combining (4.30) and (4.33), we find

$$E = \int dx dy [\frac{1}{2e^2}(B + e^2(|\phi|^2 - v^2))^2 + |(D_x - iD_y)\phi|^2 + Bv^2] \quad (4.34)$$

$$\geq \int dx dy Bv^2 = 2\pi n v^2 \quad (4.35)$$

where the inequality is saturated if and only if

$$B + e^2(|\phi|^2 - v^2) = 0, \quad (D_x - iD_y)\phi = 0. \quad (4.36)$$

In a similar manner, we can show

$$E = \int dx dy [\frac{1}{2e^2}(B - e^2(|\phi|^2 - v^2))^2 + |(D_x + iD_y)\phi|^2 - Bv^2] \geq - \int dx dy Bv^2 = -2\pi n v^2 \quad (4.37)$$

where the inequality is saturated if and only if

$$B - e^2(|\phi|^2 - v^2) = 0, \quad (D_x + iD_y)\phi = 0. \quad (4.38)$$

We conclude that $E \geq 2\pi|n|v^2$, and therefore the minimal energy configuration for a given n can be found by solving (4.36) when $n > 0$ and (4.38) when $n < 0$.

Let us study these minimal energy configuration further. Let us assume $n > 0$. We use the same trick as we used in Sec. 3.3: recalling $B = \partial_x A_y - \partial_y A_x$, we introduce ρ by demanding $A_x = -\partial_y \rho$ and $A_y = +\partial_x \rho$, meaning that $B = (\partial_x^2 + \partial_y^2)\rho$. We now define $f(x, y)$ by

$$\phi(x, y) = f(x, y) \exp \rho(x, y). \quad (4.39)$$

Then

$$(D_x - iD_y)\phi(x, y) = 0 \iff (\partial_x - i\partial_y)f(x, y) = 0. \quad (4.40)$$

which means that $f(x, y)$ is a holomorphic function of $w = x + iy$.

We also know that the winding number of ϕ at asymptotic infinity is n , meaning that $f(x, y)$ is a polynomial of degree n . The general solution is then

$$\phi(x, y) = (\exp \rho(x, y)) \prod_{i=1}^n (w - w_i) \quad (4.41)$$

where w_i can be thought of as the positions of the cores of the vortices. We then have to solve

$$B + e^2(|\phi|^2 - v^2) = 0 \quad (4.42)$$

which becomes

$$(\partial_x^2 + \partial_y^2)\rho + e^2((\exp 2\rho) \prod_i |w - w_i|^2 - v^2) = 0. \quad (4.43)$$

Here, the boundary condition on ρ is that the second term vanishes at infinity. It is known that this differential equation for ρ has a unique solution for arbitrary choices of $w_{1,\dots,n}$.

This means that for any choice of $w_{1,\dots,n}$, we can find the n -vortex solution with exactly the same energy $E = 2\pi n v^2$. This means that the forces between two minimal vortices are exactly zero and they can be superimposed. The analysis so far is completely classical, and the quantum corrections destroys the exact degeneracy. There are versions of the Abelian Higgs model with fermions which has supersymmetry. There, this exact degeneracy is known to survive quantum corrections.

Exercise 4.6. Look for the literature and learn the proof that there is a unique solution to (4.43) for any $w_{1,\dots,n}$.

Exercise 4.7. Learn how to determine ρ numerically when $w_{1,\dots,n}$ are given as the input. Plot the resulting solutions for various choices and make animations.

Exercise 4.8. No analytic solutions to ρ is known on a flat space, but it is known that the equation (4.43) simplifies on the Poincaré disk and admits simple analytic solutions. Look for the literature and study it.

4.6 Higgs mechanism in non-Abelian gauge theory

4.6.1 Scalar in the doublet

Let us now consider the system of $SU(2)$ gauge field A_μ^a and a complex scalar field ϕ^a in the doublet. We consider the action

$$S = - \int d^4x \left(\frac{1}{2g^2} \text{tr} F_{\mu\nu} F^{\mu\nu} + (D_\mu \phi)^\dagger D^\mu \phi + V(\phi) \right) \quad (4.44)$$

where

$$D_\mu \phi = (\partial_\mu + iA_\mu)\phi \quad (4.45)$$

as before and

$$V(\phi) = m^2 \phi^\dagger \phi + \frac{\lambda}{2} (\phi^\dagger \phi)^2. \quad (4.46)$$

When $m^2 > 0$, the lowest energy is achieved at $\langle \phi \rangle = 0$ and we simply have the massless SU(2) gauge field and the massive scalar field ϕ . When $m^2 < 0$, the value $|\langle \phi \rangle| = |m|/\sqrt{\lambda} =: v$ gives the lowest energy as before. For definiteness, let us choose

$$\langle \phi^a \rangle = \begin{pmatrix} v \\ 0 \end{pmatrix} \quad (4.47)$$

An infinitesimal gauge variation by

$$g = e^{i\epsilon X}, \quad X = \begin{pmatrix} X^3 & X^1 + iX^2 \\ X^1 - iX^2 & -X^3 \end{pmatrix} \quad (4.48)$$

gives

$$g \begin{pmatrix} v \\ 0 \end{pmatrix} = \begin{pmatrix} v \\ 0 \end{pmatrix} + \epsilon v \begin{pmatrix} iX^3 \\ X^1 + iX^2 \end{pmatrix}. \quad (4.49)$$

This means that the fluctuation of ϕ along the imaginary part of ϕ^1 and along both the real part and the imaginary part of ϕ^2 are gauge modes and unphysical. The only physical mode is along the real part of ϕ^1 , which we decide to call ρ . The mass can be calculated as before, and finds

$$m_\rho^2 = 2|m^2|. \quad (4.50)$$

Let us now study the gauge fields. We write

$$A_\mu = \frac{1}{2} \begin{pmatrix} A^3 & A^+ \\ A^- & -A^3 \end{pmatrix}, \quad A^\pm = A^1 \pm iA^2. \quad (4.51)$$

We expand the action and find that the part $|D_\mu \phi|^2$ provides a term quadratic in A_μ :

$$\frac{1}{2g^2} F_{\mu\nu} F^{\mu\nu} + (D_\mu \phi)^\dagger D^\mu \phi \supset \frac{1}{4g^2} F_{\mu\nu}^3 F^{3,\mu\nu} + \frac{1}{4} A_\mu^3 A^{3,\mu} |v|^2 + \frac{1}{4g^2} F_{\mu\nu}^+ F^{-,\mu\nu} + \frac{1}{4} A_\mu^+ A^{-,\mu} |v|^2. \quad (4.52)$$

Three components A^3, A^\pm of the gauge fields all have the same mass

$$m_A^2 = \frac{1}{2} g^2 v^2 = \frac{g^2 |m|^2}{2\lambda}. \quad (4.53)$$

Let us count the degrees of freedom:

- When $m^2 > 0$, three massless vectors, each with two polarizations 3×2 , together with two complex scalars, each with two real components 2×2 .
- When $m^2 < 0$, three massive vectors, each with three polarizations 3×3 , together with one real Higgs scalar ρ .

We see $3 \times 2 + 2 \times 2 = 3 \times 3 + 1$. Three vector bosons ate three scalars out of four scalars and became massive.

4.6.2 Scalar in the triplet

We next consider the case when the scalar φ is in the real triplet representation of $SU(2)$. We represent φ as a traceless hermitean 2×2 matrix

$$\varphi = \begin{pmatrix} \varphi^3 & \varphi^+ \\ \varphi^- & -\varphi^3 \end{pmatrix}, \quad \varphi^\pm = \varphi^1 \pm i\varphi^2 \quad (4.54)$$

with the gauge transformation

$${}^g\varphi = g\varphi g^{-1}. \quad (4.55)$$

The covariant derivative can then be found from the general rule to be

$$D_\mu\varphi = \partial_\mu\varphi + i[A_\mu, \varphi]. \quad (4.56)$$

We now consider the action

$$S = \int d^4x \left(\frac{1}{2g^2} \text{tr} F_{\mu\nu} F^{\mu\nu} + \frac{1}{4} \text{tr} (D_\mu\varphi)^\dagger D^\mu\varphi \right) + V(\varphi) \quad (4.57)$$

where

$$V(\varphi) = \frac{m^2}{2} \left(\frac{1}{2} \text{tr} \varphi^\dagger\varphi \right) + \frac{\lambda}{4} \left(\frac{1}{2} \text{tr} \varphi^\dagger\varphi \right)^2. \quad (4.58)$$

Let us consider the case $m^2 < 0$. The potential as always is minimal at

$$\frac{1}{2} \text{tr} \varphi^\dagger\varphi = (\varphi^1)^2 + (\varphi^2)^2 + (\varphi^3)^2 = \frac{m^2}{\lambda} =: v^2. \quad (4.59)$$

For definiteness we take

$$\langle\varphi\rangle = \begin{pmatrix} v & 0 \\ 0 & -v \end{pmatrix}. \quad (4.60)$$

We perform an infinitesimal gauge transformation by (4.48) and find

$${}^g \begin{pmatrix} v & 0 \\ 0 & -v \end{pmatrix} = \begin{pmatrix} v & 0 \\ 0 & -v \end{pmatrix} + 2\epsilon v \begin{pmatrix} 0 & X^2 - iX^1 \\ X^2 + iX^1 & 0 \end{pmatrix}. \quad (4.61)$$

This means that φ^1 and φ^2 are now gauge modes and unphysical. The mode $\rho = \delta\varphi^3$ is physical and has the mass

$$m_\rho^2 = 2|m^2| \quad (4.62)$$

as always.

Let us find the masses of the gauge fields:

$$\frac{1}{2g^2} \text{tr} F_{\mu\nu} F^{\mu\nu} + \frac{1}{2} \text{tr} (D_\mu\varphi)^\dagger D^\mu\varphi \supset \frac{1}{4g^2} F_{\mu\nu}^3 F^{3,\mu\nu} + \frac{1}{4g^2} F_{\mu\nu}^+ F^{-,\mu\nu} + \frac{1}{2} A_\mu^+ A^{-,\mu} |v|^2. \quad (4.63)$$

This means that

$$m_{A^3}^2 = 0, \quad m_{A^\pm}^2 = 2g^2|v|^2 = g^2|m^2|/\lambda. \quad (4.64)$$

Here we can see the important fact: the gauge transformation of the vev $\langle\varphi\rangle$ by the infinitesimal generator X^3 was zero, while the transformation by $X^{1,2}$ were nonzero. Correspondingly, the gauge field A^3 remained massless, while $A^{1,2}$ became massive.

The gauge transformations which keeps the vev are called *unbroken* and those which change the vev are called *broken*. We can summarize our finding by saying that unbroken gauge bosons remain massless, while broken gauge bosons become massive.

Before moving on, let us count the degrees of freedom:

- When $m^2 > 0$, we had three massless vectors, with 3×2 degrees of freedom. We also had three massive scalars.
- When $m^2 < 0$, we have one massless vector and two massive vectors, $2 + 3 \times 2$. We have one massive Higgs boson.

Again we find $3 \times 2 + 3 = 2 + 3 \times 2 + 1$.

Exercise 4.9. Study the Higgs mechanism of $SU(2)$ gauge field by scalar fields in other representations.

4.7 Monopole solution

Let us continue our study of the Higgs mechanism by the triplet scalar. We took

$$\langle\varphi\rangle = \begin{pmatrix} v & 0 \\ 0 & -v \end{pmatrix} \quad (4.65)$$

as constant across spacetime above. We can try to vary it. We do not like to have too much energy, so we would like to keep

$$\langle\varphi^1\rangle^2 + \langle\varphi^2\rangle^2 + \langle\varphi^3\rangle^2 = v^2, \quad (4.66)$$

at least almost all of the regions of the space.

For simplicity let us consider a ‘spherically symmetric’ situation where

$$\langle\varphi^{1,2,3}\rangle \sim v \frac{(x, y, z)}{r} \quad (4.67)$$

when $r = \sqrt{x^2 + y^2 + z^2}$ is big. Note that the superscript of φ^i is a index for the gauge symmetry, while on the right hand side we are using the spacetime index.

This configuration is topologically nontrivial in the sense that the sphere S^2 at spatial infinity to the field space S^2 defined by (4.66);⁶ we can compare this situation to the analysis of the vortex in Sec. 4.4, where the ‘sphere’ S^1 at spatial infinity around the vortex was mapped to the sphere S^1 of the field space defined by $|\phi|^2 = v^2$. In both cases, the configurations are topologically

⁶The sphere in \mathbb{R}^{d+1} is denoted by S^d .

distinguished by its winding number. The hedgehog solution (4.67) has the winding number 1, whereas the constant solution (4.65) has the winding number 0.

Let us show that the configuration (4.67) is a magnetic monopole. To see this, let us first parameterize the sphere at the spatial infinity by the spherical coordinate so that

$$\frac{(x, y, z)}{r} = (\cos \theta, \sin \theta \cos \phi, \sin \theta \sin \phi) \quad (4.68)$$

which corresponds to the vev

$$\langle \varphi \rangle = \begin{pmatrix} \cos \theta & \sin \theta e^{i\phi} \\ \sin \theta e^{-i\phi} & -\cos \theta \end{pmatrix} v. \quad (4.69)$$

Consider the gauge transformation

$$g_N(\theta, \phi) = \exp\left(i\frac{\theta}{2} \begin{pmatrix} 0 & ie^{i\phi} \\ -ie^{-i\phi} & 0 \end{pmatrix}\right). \quad (4.70)$$

This gauge transformation brings the vev to (4.65).

Exercise 4.10. Check this, i.e. check that $g_N \langle \varphi \rangle = g_N \langle \varphi \rangle g_N^{-1} = \begin{pmatrix} v & 0 \\ 0 & -v \end{pmatrix}$.

But this gauge transformation is not single valued at the south pole, where $\theta = \pi$, for which we have

$$g_N(\pi, \phi) = \begin{pmatrix} 0 & -e^{i\phi} \\ e^{i\phi} & 0 \end{pmatrix} \quad (4.71)$$

Instead, around very close to the south pole, we need to choose a single-valued gauge transformation

$$g_S(\pi, \phi) = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}. \quad (4.72)$$

Between the two transformations, we find

$$g_N(\pi, \phi) = \begin{pmatrix} e^{i\phi} & 0 \\ 0 & e^{-i\phi} \end{pmatrix} g_S(\pi, \phi), \quad (4.73)$$

where we note that the gauge transformations between the two patches,

$$\begin{pmatrix} e^{i\phi} & 0 \\ 0 & e^{-i\phi} \end{pmatrix}, \quad (4.74)$$

is in the U(1) subgroup unbroken by the vev $\begin{pmatrix} v & 0 \\ 0 & -v \end{pmatrix}$. Recalling the arguments in Sec. 3.1, the existence of this gauge transformation of the winding number 1 means that the monopole flux in the unbroken U(1) gauge field is

$$\int d\vec{n} \cdot \vec{B} = \int d\theta d\phi F_{\theta\phi} = 2\pi \cdot 1. \quad (4.75)$$

In general, the winding number of the scalar field at spatial infinity, $S^2 \xrightarrow{\varphi} S^2$ translates to the monopole number.

Exercise 4.11. Derive this result by a topological argument.

Analogously to the case of the vortex, $\langle \varphi \rangle$ needs to vanish at the core of the vortex. The profile of the scalar fields can be found by solving the equations of motion derived from the action. This solution is called as the 't Hooft-Polyakov monopole, after the discoverers.

Exercise 4.12. Study the literature and learn how to determine the profile of the monopole.

There is also the version of the Bogomolny trick for the monopole.

$$E \geq \int d^3x \frac{1}{g^2} \text{tr} B_i B_i + \frac{1}{4} (D_i \varphi) D_i \varphi \quad (4.76)$$

$$= \int d^3x \text{tr} \left(\frac{B_i}{g} \pm \frac{1}{2} D_i \varphi \right)^2 \mp \text{tr} \frac{1}{g} B_i D_i \varphi \quad (4.77)$$

$$\geq \mp \frac{1}{g} \int d^3x \text{tr} D_i (B_i \varphi) \quad (4.78)$$

$$= \mp \frac{1}{g} \int dS_i \text{tr} B_i \varphi \quad (4.79)$$

$$= \mp \frac{v}{g} \int d\vec{S} \cdot \vec{B} = \mp \frac{v}{g} 2\pi n = \frac{v}{g} 2\pi |n|. \quad (4.80)$$

In the first inequality, we dropped the contribution from $V(\varphi)$. In the third line, we used the fact $D_i B_i = 0$ and added $\text{tr}(D_i B_i) \varphi_i$; this can be proved easily since $B_i = \epsilon_{ijk} [D_j, D_k]/2$ and therefore

$$D_i B_i = \epsilon_{ijk} [D_i, [D_j, D_k]]/2 = 0. \quad (4.81)$$

In the fourth line, we convert the volume integral to the surface integral at spatial infinity. To go to the fifth line, we use the fact that $\text{tr} B_i \varphi$ project the non-Abelian gauge field to the unbroken $U(1)$ direction; recall our normalization is such that

$$B_i = \frac{1}{2} \begin{pmatrix} B_i^{\text{unbroken}} & 0 \\ 0 & -B_i^{\text{unbroken}} \end{pmatrix}, \quad \langle \varphi \rangle = \begin{pmatrix} v & 0 \\ 0 & -v \end{pmatrix}. \quad (4.82)$$

Finally, we note that the second inequality is attained only when

$$\frac{B_i}{g} = \mp \frac{1}{2} D_i \varphi. \quad (4.83)$$

This is also called the Bogomolny equation. For a general potential $V(\varphi)$, this equation does not help much, since we discarded $V(\varphi)$ at the first inequality. But for the special case where $V(\varphi) \equiv 0$, the solutions of the Bogomolny equation give the lowest energy monopole configurations for a given monopole number n .

4.8 Higgs mechanism in the Standard model

4.8.1 Gauge boson masses

In the standard model, we have $U(1) \times SU(2)$ gauge fields and a Higgs field ϕ with the covariant derivative

$$D_\mu \phi = (\partial_\mu + \frac{i}{2} A_\mu^{U(1)} + i A_\mu^{SU(2)}) \phi. \quad (4.84)$$

As a matrix, the gauge field part has the form

$$\left(\frac{1}{2} A_\mu^{U(1)} + A_\mu^{SU(2)}\right) \phi = \frac{1}{2} \begin{pmatrix} A_\mu^{U(1)} + A_\mu^3 & A_\mu^+ \\ A_\mu^- & A_\mu^{U(1)} - A_\mu^3 \end{pmatrix} \begin{pmatrix} \phi^1 \\ \phi^2 \end{pmatrix}. \quad (4.85)$$

We now assume that our potential $V(\phi)$ is such that it has the minimum at $|\langle \phi \rangle| = v = m/\sqrt{\lambda}$. We pick the convention that

$$\langle \phi \rangle = \begin{pmatrix} 0 \\ v \end{pmatrix}. \quad (4.86)$$

The gauge field masses can be found from

$$\begin{aligned} & \frac{1}{4g_1^2} F_{\mu\nu}^{U(1)} F^{U(1),\mu\nu} + \frac{1}{2g_2^2} \text{tr} F_{\mu\nu}^{SU(2)} F^{SU(2),\mu\nu} + (D_\mu \phi)^\dagger D^\mu \phi \supset \\ & \frac{1}{4g_1^2} F_{\mu\nu}^{U(1)} F^{U(1),\mu\nu} + \frac{1}{4g_2^2} F_{\mu\nu}^3 F^{3,\mu\nu} + \frac{1}{4g_2^2} F_{\mu\nu}^+ F^{-,\mu\nu} \\ & + \frac{1}{4} (A_\mu^{U(1)} - A_\mu^3) (A^{U(1),\mu} - A^{3,\mu}) |v|^2 + \frac{1}{4} A_\mu^+ A^{-,\mu} |v|^2 \end{aligned} \quad (4.87)$$

The component A^\pm are usually called the W-bosons and denoted by W^\pm . They have masses

$$m_W^2 = \frac{1}{2} g_2^2 |v|^2. \quad (4.88)$$

We now define

$$\begin{pmatrix} A_\mu^{U(1)} \\ A_\mu^3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} A_\mu^{\text{em}} + \begin{pmatrix} g_1/g_2 \\ -g_2/g_1 \end{pmatrix} Z_\mu. \quad (4.89)$$

Then we have

$$\frac{1}{4g_1^2} F_{\mu\nu}^{U(1)} F^{U(1),\mu\nu} + \frac{1}{4g_2^2} F_{\mu\nu}^3 F^{3,\mu\nu} + \frac{1}{4} (A_\mu^{U(1)} - A_\mu^3) (A^{U(1),\mu} - A^{3,\mu}) |v|^2 \quad (4.90)$$

$$= \frac{1}{4} \left(\frac{1}{g_1^2} + \frac{1}{g_2^2} \right) F_{\mu\nu}^{\text{em}} F^{\text{em},\mu\nu} + \frac{1}{4} \left(\frac{1}{g_1^2} + \frac{1}{g_2^2} \right) Z_{\mu\nu} Z^{\mu\nu} + \frac{1}{4} \left(\frac{g_1}{g_2} + \frac{g_2}{g_1} \right)^2 Z_\mu Z^\mu |v|^2 \quad (4.91)$$

$$= \frac{1}{4e^2} F_{\mu\nu}^{\text{em}} F^{\text{em},\mu\nu} + \frac{1}{4} \frac{g_1^2 + g_2^2}{g_1^2 g_2^2} Z_{\mu\nu} Z^{\mu\nu} + \frac{1}{4} \frac{(g_1^2 + g_2^2)^2}{g_1^2 g_2^2} Z_\mu Z^\mu |v|^2. \quad (4.92)$$

We see that A_μ^{em} remains massless; this is the electromagnetic $U(1)$ we actually observe, and the coupling e is given by

$$e = \frac{g_1 g_2}{\sqrt{g_1^2 + g_2^2}}. \quad (4.93)$$

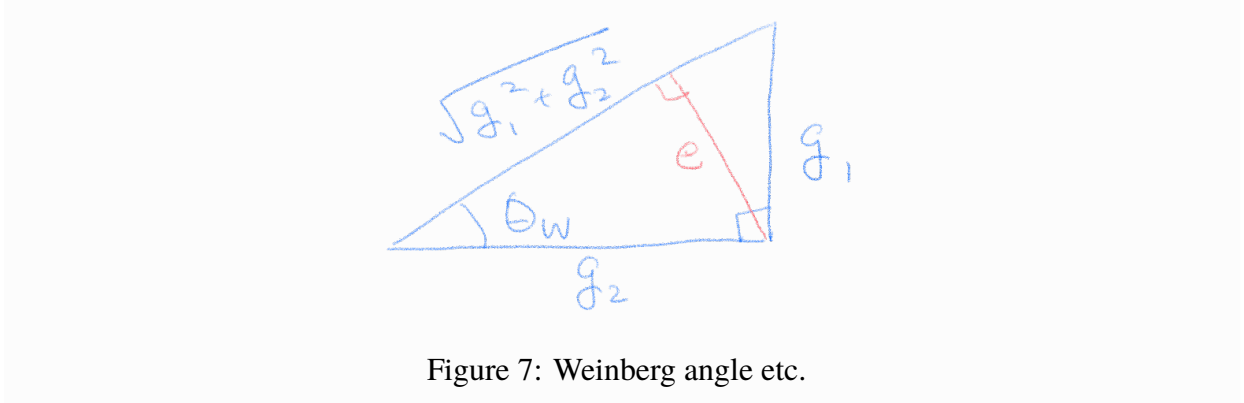


Figure 7: Weinberg angle etc.

The mass of the Z-boson is given by

$$m_Z^2 = \frac{1}{2}(g_1^2 + g_2^2)|v|^2. \quad (4.94)$$

We see that

$$\frac{m_W}{m_Z} = \cos \theta_W, \quad \tan \theta_W = \frac{g_1}{g_2}. \quad (4.95)$$

This angle θ_W is known as the Weinberg angle. See Fig. 7.

4.8.2 Fermion masses

Recall from Sec. 2.7.4 that the Standard Model Lagrangian did not have a mass term for the quarks. What we had was the Yukawa interaction between fermions and the Higgs:

$$(\text{quark Yukawa}) = Y_{ij}^{\text{up}} \epsilon_{uv} \delta_{a\bar{a}} \phi^u (Q_L)^{iva\alpha} \bar{u}_{R\alpha}^{j\bar{a}} + Y_{ij}^{\text{down}} \delta_{\bar{u}v} \delta_{a\bar{a}} \bar{\phi}^{\bar{u}} (Q_L)^{iva\alpha} \bar{d}_{R\alpha}^{j\bar{a}} + c.c. \quad (4.96)$$

We now use

$$Q_L^{iva\alpha} = \begin{pmatrix} u_L^{ia\alpha} \\ d_L^{ia\alpha} \end{pmatrix}. \quad (4.97)$$

and plug in the Higgs vev to this Yukawa interaction. We find

$$(\text{quark Yukawa}) \supset v Y_{ij}^{\text{up}} \delta_{a\bar{a}} \phi^u (u_L)^{ia\alpha} \bar{u}_{R\alpha}^{j\bar{a}} + \bar{v} Y_{ij}^{\text{down}} \delta_{a\bar{a}} (d_L)^{ia\alpha} \bar{d}_{R\alpha}^{j\bar{a}} + c.c. \quad (4.98)$$

This means that they become the mass terms for the quarks.

Before proceeding, we note that the covariant derivative of Q_L contained the combination

$$\partial_\mu + i \frac{1}{6} A_\mu^{\text{U}(1)} + i A_\mu^{\text{SU}(2)} = \partial_\mu + i \begin{pmatrix} \frac{1}{6} A_\mu^{\text{U}(1)} + \frac{1}{2} A_\mu^3 & & \frac{1}{2} A_\mu^+ \\ & \frac{1}{2} A_\mu^- & \\ & & \frac{1}{6} A_\mu^{\text{U}(1)} - \frac{1}{2} A_\mu^3 \end{pmatrix} \quad (4.99)$$

acting on (4.97). Using (4.89), the coupling to A_μ^{em} can be found by simply setting $A_\mu^{\text{U}(1)} = A_\mu^3 = A_\mu^{\text{em}}$. This means that u_L and d_L have the electromagnetic U(1) charges $1/6 + 1/2 = 2/3$ and $1/6 - 1/2 = -1/3$, respectively. Similarly, the electromagnetic U(1) charges of \bar{u}_R and \bar{d}_R can be found to be $-2/3$ and $+1/3$. Indeed the mass terms are invariant under the electromagnetic U(1) gauge symmetry.

4.8.3 Number of physical parameters in the Yukawa couplings

Let us count the number of physical parameters in the Yukawa couplings of the quarks. (The following discussion can be done without referring to the Higgs effect, so this part could have been put at the end of Sec. 2.7.4.) For this purpose, it is useful to write the quark Yukawa by emphasizing a slightly different aspect:

$$(\text{quark Yukawa}) = Y_i^x \phi(Q_L)^i (\overline{u_R})_x + \tilde{Y}_i^s \bar{\phi}(Q_L)^i (\overline{d_R})_s + c.c. \quad (4.100)$$

where we dropped the gauge indices u, a and the spinor index α because they do not play any role below, and we used distinct symbols $i, x, s = 1, \dots, N$ for the label of the generation (which is the real world is given by $N = 3$).

We now note that it is our choice to redefine the symbols as follows:

$$Q_L^i \mapsto \mathcal{U}_j^i Q_L^j, \quad (\overline{u_R})_x \mapsto (\overline{u_R})_y (\mathcal{V}^{-1})_x^y, \quad (\overline{d_R})_s \mapsto (\overline{d_R})_t (\mathcal{V}^{-1})_s^t \quad (4.101)$$

and at the same time

$$Y_i^x \mapsto \mathcal{V}_y^x Y_j^y (\mathcal{U}^{-1})_i^j, \quad \tilde{Y}_i^s \mapsto \mathcal{W}_t^s \tilde{Y}_j^t (\mathcal{U}^{-1})_i^j \quad (4.102)$$

where $\mathcal{U}, \mathcal{V}, \mathcal{W}$ are three unitary matrices which are independent of the spacetime position.

This means that any observable physical quantity is a function of Y, \tilde{Y} and their complex conjugates, which are invariant under the transformation (4.102).⁷

The invariance under \mathcal{V} and \mathcal{W} are easy to deal with: we simply form

$$X_j^i := (Y^\dagger)_x^i Y_j^x, \quad \tilde{X}_j^i := (\tilde{Y}^\dagger)_s^i \tilde{Y}_j^s \quad (4.103)$$

which are both hermitean and have the transformation

$$X \mapsto \mathcal{U} X \mathcal{U}^{-1}, \quad \tilde{X} \mapsto \mathcal{U} \tilde{X} \mathcal{U}^{-1}. \quad (4.104)$$

An $N \times N$ Hermitean matrix contains N^2 real parameters. In total, X and \tilde{X} have $2N^2$ real parameters. There is an action of \mathcal{U} , containing N^2 real parameters, but \mathcal{U} and $c\mathcal{U}$ act in the same way on X and \tilde{X} . In total we find

$$2N^2 - (N^2 - 1) = N^2 + 1 \quad (4.105)$$

parameters.

Note that the eigenvalues of X, \tilde{X} (times v^2) are the masses squared of the up-type quarks and the down-type quarks, respectively. They comprise $2N$ real parameters.

Therefore, there are

$$N^2 + 1 - 2N = (N - 1)^2 \quad (4.106)$$

additional real parameters, which describe how the quarks are mixed under the weak interaction. For $N = 2$ there is a single such parameter, known as the Cabbibo angle. For $N = 3$, there are $2^2 = 4$ parameters.

⁷In mathematics this is an example of something called the quiver representation theory.

4.8.4 CP (non-)invariance

A CP transformation on a fermion acts in the following way: $\psi_\alpha^{\text{new}}(t, x, y, z) = \overline{\psi_\alpha(t, -x, -y, -z)}$, which is compatible with the Lorentz transformation.

Exercise 4.13. Check this.

One finds that performing the CP transformation on the fermions is equivalent to performing the change

$$Y_i^{\text{new}x} := \overline{Y_i^x}, \quad \tilde{Y}_i^{\text{new}s} := \overline{\tilde{Y}_i^s}. \quad (4.107)$$

We say that our theory is CP invariant when Y and \tilde{Y} are invariant under (4.107), i.e. they are real matrices, possibly after a suitable transformation (4.102) using \mathcal{U} , \mathcal{V} and \mathcal{W} . The discussion in the last subsection means that this condition is equivalent to the condition that X and \tilde{X} can be made real matrices by a suitable \mathcal{U} .

Let us count the number of parameters in CP invariant theories. A real symmetric matrix X contains $N(N + 1)/2$ parameters; and therefore X and \tilde{X} together have $N(N + 1)$ parameters. We can still perform the transformation (4.104) with a real orthogonal matrix \mathcal{U} , which contains $N(N - 1)/2$ parameters. In the end we find

$$N(N + 1) - N(N - 1)/2 = N(N + 3)/2 \quad (4.108)$$

parameters. Among them, we have $2N$ mass parameters. Therefore there are

$$N(N + 3)/2 - 2N = N(N - 1)/2 \quad (4.109)$$

mixing angles. Comparing with (4.106), we find that the non-CP-invariant theories have

$$(N - 1)^2 - N(N - 1)/2 = (N - 1)(N - 2)/2 \quad (4.110)$$

real parameters in addition to the CP-invariant cases.

When $N = 2$, this is zero. This means that with two generations, the quark sector is automatically CP invariant. When $N = 3$, this is one. This means that with three generations, the quark sector contains a single real parameter which controls its CP non-invariance.⁸

One consequence of the CP invariance is the following. Recall that any physical quantity is a function of X and \tilde{X} invariant under the transformation (4.104). Any function can be approximated

⁸This simple analysis gave Kobayashi and Maskawa a Nobel Prize! I always consider a 4-way distinction concerning scientific problems:

	interesting	not interesting
simple	A	B
difficult	C	D

I like A the most; the result of Kobayashi-Maskawa clearly belongs to this category. D is the worst. I am not sure which of B and C I prefer.

by a polynomial. Any monomial made out of X and \tilde{X} invariant under the transformation (4.104) has the form

$$\text{tr}(\text{a sequence of } X \text{ and } \tilde{X},) \quad (4.111)$$

for example

$$\text{tr } XX\tilde{X}X\tilde{X}\tilde{X}. \quad (4.112)$$

When the theory is CP invariant, X and \tilde{X} can be made simultaneously real, and such traces are automatically real. When the theory is not CP invariant, there are in general non-real combinations among such functions.

We saw above that the system is always CP-invariant when $N = 2$. This means that for arbitrary pair of two 2×2 Hermitean matrices X and \tilde{X} , the trace of an arbitrary sequence of X and \tilde{X} , for example (4.112), is automatically real.

Exercise 4.14. Prove this funny fact about two 2×2 Hermitean matrices.

5 Renormalization group

In the last section we only talked about purely classical phenomena. Let us study some quantum effects.

5.1 Scalar ϕ^4

Let us start from a theory of N real scalars $\phi_{i=1,\dots,N}$, with the action

$$S_E = \int d^4x \left[\frac{1}{2} \partial_\mu \phi_i \partial^\mu \phi_i + \frac{1}{24} \lambda_L (\phi_i \phi_i)^2 \right] \quad (5.1)$$

where we work in the Wick-rotated Euclidean version; the repeated indices are summed. We also put the subscript L to the parameter λ that it is a parameter in the Lagrangian which is not directly the quantity appearing in the correlation functions. Before proceeding, we note that ϕ has dimension 1 since d^4x has dimension -4 and ∂ has dimension 1. Therefore the coupling λ_L is dimensionless.

At tree level, we find

$$\langle \phi_i(p) \phi_i(q) \phi_i(r) \phi_i(s) \rangle \supset (2\pi)^4 \delta(p + q + r + s) \frac{1}{p^2} \frac{1}{q^2} \frac{1}{r^2} \frac{1}{s^2} (-\lambda_L) \quad (5.2)$$

where we *do not* sum over the repeated indices i . and we only kept the term proportional to $1/p^2 q^2 r^2 s^2$. See Fig. 8.0) for the diagram.

There are a few one-loop diagram. The two diagrams Fig. 8.1) and 1') contribute by

$$(2\pi)^4 \delta(p + q + r + s) \frac{1}{p^2} \frac{1}{q^2} \frac{1}{r^2} \frac{1}{s^2} (-\lambda_L)^2 \left(\frac{1}{2} + \frac{N-1}{18} \right) \int d^4\ell \frac{1}{\ell^2 (\ell + p + q)^2}. \quad (5.3)$$

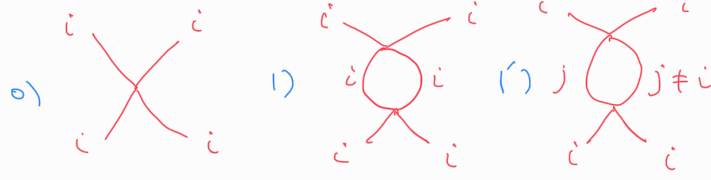


Figure 8: Diagrams for the renormalization of the ϕ^4 theory

These two diagrams group momenta as $(pq)(rs)$. There are also diagrams where we group momenta as $(pr)(qs)$ and $(ps)(qr)$.

Let us concentrate on the integral $\int d^4\ell/(\ell^2(\ell + P)^2)$. This itself is logarithmically divergent, but its derivative w.r.t. P^2 is not. Let us perform some standard manipulations

$$\begin{aligned} I(P^2) &:= \int \frac{d^4\ell}{(2\pi)^4} \frac{1}{\ell^2(\ell + P)^2} \\ &= \int_0^1 dx \int \frac{d^4\ell}{(2\pi)^4} \frac{1}{(\ell^2 + x(1-x)P^2)^2} = \int_0^1 dx \frac{2\pi^2}{(2\pi)^4} \int_0^\infty \ell^3 d\ell \frac{1}{(\ell^2 + x(1-x)P^2)^2}. \end{aligned} \quad (5.4)$$

We now formally differentiate the integrand with respect to P^2 . After a short computation one finds

$$P^2 \frac{\partial}{\partial P^2} I = -\frac{1}{16\pi^2}. \quad (5.5)$$

Exercise 5.1. Perform this computation.

We now consider the dependence of $\langle \phi_i(p)\phi_i(q)\phi_i(r)\phi_i(s) \rangle$ when we keep the ratio $p : q : r : s$ fixed and change the scale $(p+q)^2 \sim (p+r)^2 \sim (p+s)^2 \sim P^2$. We denote by $\lambda(P^2)$ the coefficient appearing in

$$\langle \phi_i(p)\phi_i(q)\phi_i(r)\phi_i(s) \rangle \supset (2\pi)^4 \delta(p+q+r+s) \frac{1}{p^2} \frac{1}{q^2} \frac{1}{r^2} \frac{1}{s^2} (-\lambda(P^2)). \quad (5.6)$$

The computations above show that

$$\frac{\partial}{\partial \log P} \lambda(P^2) = \frac{N+8}{3} \lambda_L^2 + O(\lambda_L^3). \quad (5.7)$$

We now note $\lambda(P^2) = \lambda_L + O(\lambda_L^2)$, and therefore we write this as

$$\frac{\partial}{\partial \log P} \lambda(P^2) = \frac{N+8}{3} \frac{1}{16\pi^2} \lambda(P^2)^2 + O(\lambda(P^2)^3). \quad (5.8)$$

Note that λ_L disappeared from our view. Or more simply, we have

$$\frac{\partial}{\partial \log P} \lambda = \frac{N+8}{3} \frac{1}{16\pi^2} \lambda^2 + O(\lambda^3). \quad (5.9)$$

This is an equation written solely in terms the quantities appearing directly in the correlation functions, and does not contain λ_L etc. which are formally infinite. This differential equation is known as the *renormalization group equation*, and the change in the coupling is known as its *running*.⁹

Neglecting the $O(\lambda^3)$ correction, it is easy to integrate it. One first rewrites it as

$$\frac{\partial}{\partial \log P} \frac{16\pi^2}{\lambda} = -\frac{N+8}{3} \quad (5.10)$$

which means that

$$\frac{16\pi^2}{\lambda} = \frac{N+8}{3} \log \frac{P_{\text{LP}}}{P} \quad (5.11)$$

where Λ_{LP} is an integration constant.

This means that λ grows as we raise the energy scale P , and becomes infinite when $P = P_{\text{LP}}$; this is called the *Landau pole*. Of course the $O(\lambda^3)$ correction become non-negligible well before we reach this scale; the point is that the theory goes out of the validity of the perturbation theory around this scale. In contrast, when we lower the energy scale to measure the system, λ smoothly goes to zero. This means that the self-interaction of the scalar field decreases as we measure the system at larger distances.

Another point of note is that the original Lagrangian (5.1) does not have any dimensionful parameter; but after solving the renormalization group equation one finds a dimension-ful parameter P_{LP} . This is sometimes called the *dimensional transmutation*.

One application of this simple analysis was an theoretical upper bound on the Higgs mass before it was discovered. Let us put ourselves before the Higgs discovery. Recall that the Higgs mass is $m_H^2 = 2|m|^2 = 2v^2\lambda$ and the W-boson mass is $m_W^2 = g^2v^2/2$. As the W-boson had been already discovered and the weak coupling g was also known, the vev of the Higgs field v was known. This means that determining m_H was equivalent to determining λ , around the energy scale $P^2 \sim m_W^2$. Using the result above with $N = 4$, we can then estimate the scale of the Landau pole P_{LP} . One might assume that P_{LP} should not be below the Planck scale M_{Planck} (or some other favorite scale of yours.) This puts an upper bound to λ . This type of bound is known as the *triviality bound* from historical reasons. When we require $P_{\text{LP}} > M_{\text{Planck}}$, we get the constraint $m_H \lesssim 150\text{GeV}$.

Exercise 5.2. Perform this computation.

In passing, we note that in the full Standard Model there are many more terms which contribute to the running of λ . For example, there is also a contribution of the top loop, see Fig. 9. Since this is a fermion loop, one has an extra minus sign in the contribution:

$$\frac{\partial}{\partial \log P} \lambda \sim +\lambda^2 - Y_t^4 + \dots \quad (5.12)$$

Again, Y_t was known since the top quark had already been discovered. When λ is small, then λ decreases as we raise the energy, and it can happen that λ becomes negative before too long. This

⁹People often say that the coupling constant runs, but it is not constant since it runs. I try to avoid to say the running coupling constants and try to simply say the running couplings, but it is hard to change the habit.

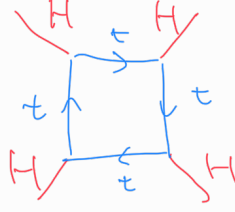


Figure 9: Top loop contribution to the running of Higgs self coupling.

makes the Higgs potential unbounded from below, and makes our present universe only metastable, or even unstable depending on the choice of the parameter. This gives a lower bound to the Higgs mass. This type of the bound was known as the *stability bound*.

5.2 Running of the gauge couplings

5.2.1 General formula

Let us move on to the case of the gauge theory. We start from the action

$$S_E = \int d^4x \left[\frac{1}{4g^2} \text{tr} F_{\mu\nu} F^{\mu\nu} + \bar{\psi} D_\mu \sigma^\mu \psi + D_\mu \phi^\dagger D^\mu \phi \right] \quad (5.13)$$

where $F_{\mu\nu}$ is a gauge field of gauge group G , ψ is a left-handed fermion in the representation R_F of G , and ϕ is a complex scalar in the representation R_S of G .

At tree level, $\langle A_\mu^a(p) A_\nu^b(-p) \rangle \sim \delta^{ab} g^2 / p^2$, where $a = 1, \dots, \dim G$ is the index for the gauge components. We have various loop corrections to this, see Fig. 10.

It is clear from the diagrams that the running is of the form

$$\frac{\partial}{\partial \log P} g^2 \sim g^4 (k_v C(\mathfrak{g}) + k_f C(R_F) + k_s C(R_S)) + O(g^6). \quad (5.14)$$

Here, $C(R)$ for an irreducible representation R is defined by

$$C(R) \delta^{ab} = \text{tr} \rho(T^a) \rho(T^b) \quad (5.15)$$

where $T^{a=1, \dots, \dim G}$ are infinitesimal generators of the Lie algebra \mathfrak{g} of G and ρ is the representation matrix in R , and for a reducible representation $C(R)$ is given by the sum $\sum_i C(R_i)$ where R_i are the irreducible components of R ; $C(R)$ are manifestly positive. Then, k_v , k_f and k_s are the numerical constants which characterize the loop computations of a vector boson, a fermion, or a scalar field. Their signs are not clear until one actually computes it.

After a long computation, one finds

$$\frac{\partial}{\partial \log P} g^2 = \frac{g^4}{8\pi^2} \left[-\frac{11}{3} C(\mathfrak{g}) + \frac{2}{3} C(R_f) + \frac{1}{3} C(R_s) \right] + O(g^6). \quad (5.16)$$

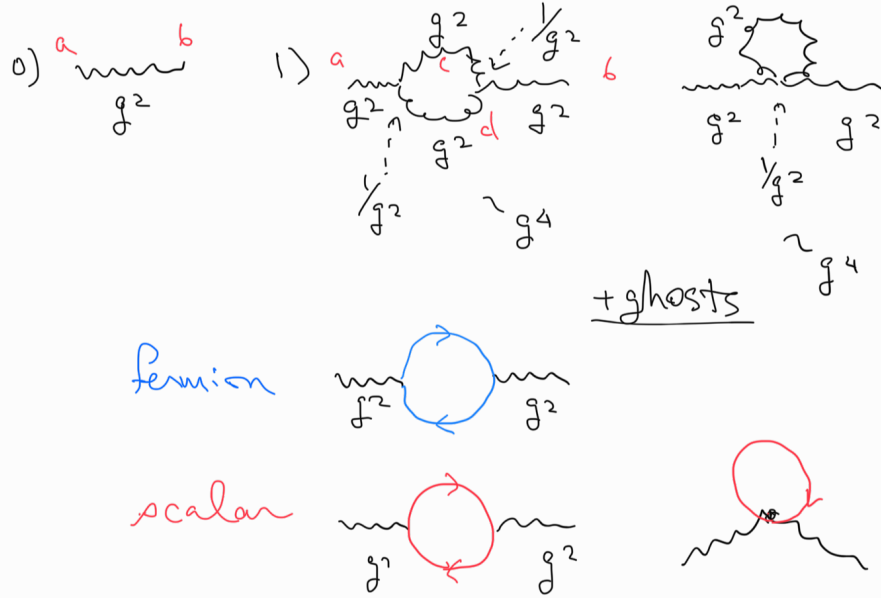


Figure 10: One loop diagrams contributing to the running of the gauge coupling.

Exercise 5.3. If you have never done this computation, you should. Consult any of the standard textbooks.

5.2.2 QED

Consider for example the $U(1)$ gauge theory with Dirac fermions Ψ_i of charge q_i , i.e. the quantum electrodynamics (QED). Each Dirac fermion is a pair of Weyl fermions ψ_L and $\bar{\psi}_R$, each of charge q_i and $-q_i$. Plugging this in to the general formula, one finds

$$\frac{\partial}{\partial \log P} e^2 = \frac{e^4}{6\pi^2} \left(\sum q_i^2 \right) + O(e^6). \quad (5.17)$$

Neglecting the higher order corrections, this is easy to integrate. It is conventional to use $\alpha := e^2/(4\pi)$, called the fine structure constant. We find

$$\frac{\partial}{\partial \log P} \frac{1}{\alpha} = -\frac{2}{3\pi} \left(\sum q_i^2 \right) + O(\alpha), \quad (5.18)$$

meaning that

$$\frac{1}{\alpha} = \frac{2}{3\pi} \left(\sum q_i^2 \right) \log \frac{P_{LP}}{P} \quad (5.19)$$

where P_{LP} is again the Landau pole. This means that the perturbation theory breaks down around the high energy scale $\sim P_{LP}$ in the massless QED, and the coupling runs to zero in the infrared. This behavior is known as the *infrared freedom*.

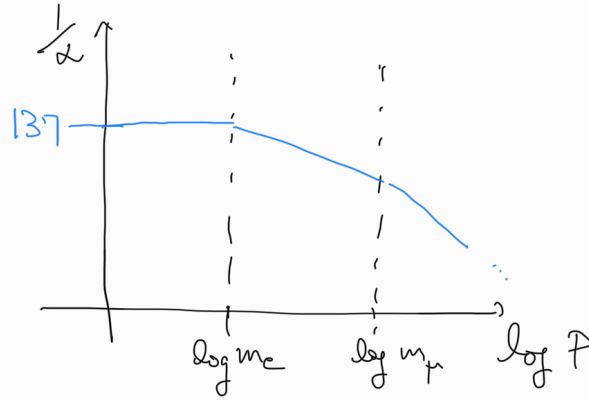


Figure 11: Running of the fine structure constant.

When the fermion has mass M , it is known that the contribution to the running disappears when $P \lll M$ and the running is essentially the same with the massless fermion when $P \ggg M$. This needs a careful justification but should be intuitively acceptable. In the real world, the lightest charged particle is the electron (which is about 500keV) and the second lightest one is the muon (which is about 150MeV). This means that the fine structure constant stops running below the scale of the electron mass; this is the value of the fine structure constant which we usually refer to, and famously has the dimensionless value $\sim 1/137$.¹⁰

5.2.3 QCD

Let us next consider the $SU(N)$ gauge theory with N_f copies of Dirac fermions in the fundamental N -dimensional representation. It is conventional to use the normalization that $C(\text{fundamental}) = 1/2$ for T^a . Then one finds that $C(\mathfrak{su}(N)) = N$.

Exercise 5.4. Check this.

Plugging into our general formula, one finds

$$\frac{\partial}{\partial \log P} g^2 = \frac{g^4}{24\pi^2} (-11N + 2N_f). \quad (5.20)$$

In particular, when $N_f = 0$ (which is called the pure Yang-Mills theory, the coupling grows in the infrared; the naive solution of the renormalization group equation is

$$\frac{1}{g^2} = \frac{11}{24\pi^2} \log \frac{P}{\Lambda_{\text{dyn}}}. \quad (5.21)$$

¹⁰There have been many crackpots who tried to come up with a formula for α . Its history up to 2003 is summarized in [Kra03], according to which Heisenberg thought $\alpha = \pi^2/(2^4 3^3)$. In 2018 it was a mildly big issue in academia that the famous mathematician M. F. Atiyah claimed to have come up with the derivation, which he announced in a public talk. He passed away shortly after, before describing the details.

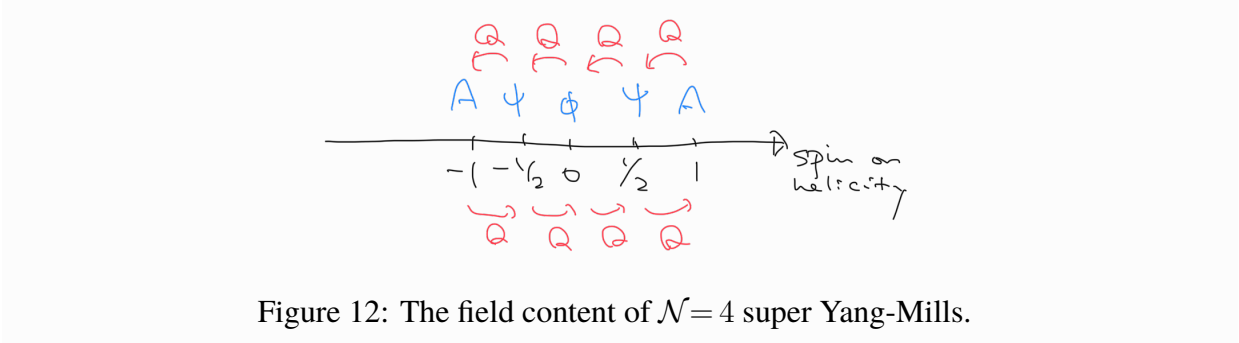


Figure 12: The field content of $\mathcal{N}=4$ super Yang-Mills.

This means that something happens at around the *dynamical scale* Λ_{dyn} . What we believe to happen is the *confinement* and the *generation of the mass gap*, to which we return in detail later. Its ultraviolet behavior in contrast is very tame: g^2 smoothly goes to zero as we raise the energy scale. This behavior is known as the *asymptotic freedom*.

The asymptotic freedom persists as long as $N_f/N_c < 11/2$. When N_f/N_c exceeds this value, the direction of the running of the coupling reverses, and the system becomes infrared free, just as in the case of massless QED. We find the Landau pole in the ultraviolet.

5.2.4 $\mathcal{N}=4$ super Yang-Mills

It so happens that our general formula gives zero when we choose four chiral fermions ψ in \mathfrak{g} and three complex scalars ϕ in \mathfrak{g} for any gauge group G , since

$$-\frac{11}{3} + 4 \times \frac{2}{3} + 3 \times \frac{1}{3} = 0. \quad (5.22)$$

It is known that by giving a suitable potential for the scalars and a suitable Yukawa coupling among fermions and scalars, the system has $\mathcal{N}=4$ supersymmetry. A supersymmetry is a symmetry which changes the spin of the particle by $1/2$; In an $\mathcal{N}=4$ supersymmetric system there are four supersymmetries, which connects all the fields A_μ , ψ_α and ϕ in the system, see Fig. 12. With $\mathcal{N}=4$ supersymmetry, there is an independent argument saying that the gauge coupling cannot run at all, not just to the leading order but to all orders. What we saw above is consistent with this curiosity.

5.3 Two-loop running and the fixed points

There are of course higher-loop contributions to the renormalization group equation. Let us first discuss some generalities.

Let us suppose that the renormalization group equation is of the form

$$\frac{\partial}{\partial \log P} g^2 = b_1 g^4 + b_2 g^6 + O(g^8), \quad (5.23)$$

where g^2 is a coupling in the theory which is supposed to be positive, such as the gauge coupling squared or the coefficient of ϕ^4 in the Higgs potential. We first need to discuss the issue of the *scheme dependence* of the renormalization group equation. This concerns the following point.

Very naively, a coupling, say g^2 , is a coefficient of a term in the Lagrangian. However, a naive loop computation leads to infinities, which needs to be regularized and renormalized. We then need to define $g^2(P^2)$ depending on the scale, in terms of a correlation function. This process involves many choices, and a totality of such choices which allows computations in the QFT under consideration is called a renormalization scheme, or simply a *scheme*. Suppose we have two schemes, and correspondingly two differently defined couplings g_{scheme1}^2 and g_{scheme2}^2 . An actually measurable quantity X is a function of the coupling, and can be computed in any scheme of your choice. Of course should satisfy

$$X = X_{\text{scheme1}}(g_{\text{scheme1}}^2) = X_{\text{scheme2}}(g_{\text{scheme2}}^2) \quad (5.24)$$

under a certain mapping $g_{\text{scheme1}} \leftrightarrow g_{\text{scheme2}}$.

Usually scientists agree among themselves on the leading term of what should be their g^2 , so the relation is usually¹¹ of the form

$$g_{\text{scheme2}}^2 = 1 \cdot g_{\text{scheme1}}^2 + c \cdot g_{\text{scheme1}}^4 + c' \cdot g_{\text{scheme1}}^6 + \dots \quad (5.25)$$

Now, suppose we computed the renormalization group equation

$$\frac{\partial}{\partial \log P} g_{\text{scheme1}}^2 = b_1 g_{\text{scheme1}}^4 + b_2 g_{\text{scheme1}}^6 + b_3 g_{\text{scheme1}}^8 + \dots \quad (5.26)$$

We can use (5.25) to translate it to the renormalization group equation for g_{scheme2} . It turns out that

$$\frac{\partial}{\partial \log P} g_{\text{scheme2}}^2 = b'_1 g_{\text{scheme2}}^4 + b'_2 g_{\text{scheme1}}^6 + b'_3 g_{\text{scheme1}}^8 + \dots \quad (5.27)$$

where

$$b'_1 = b_1, \quad b'_2 = b_2, \quad b'_3 = b_3 - b_2 c + b_1(c' - c^2) \dots \quad (5.28)$$

This means that the first two coefficients of the renormalization group equation are scheme independent.

Exercise 5.5. Confirm (5.28).

Exercise 5.6. Come up with an argument that $b_{1,2}$ are scheme independent which does not use an explicit computation.

Exercise 5.7. Does this scheme independence apply when there are more than one couplings?

¹¹This is not always the case. If you work on supersymmetry, the often-used holomorphic scheme does not satisfy this.

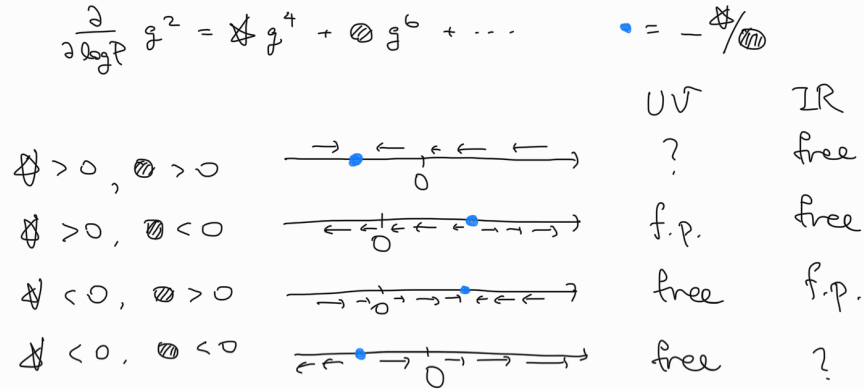


Figure 13: Schematic form of the renormalization group flow.

Let us now come back to the equation

$$\frac{\partial}{\partial \log P} g^2 = b_1 g^4 + b_2 g^6 + O(g^8). \quad (5.29)$$

We dropped the subscript specifying the scheme again. A rough drawing of the renormalization group flow of g^2 is given in Fig. 13. Note that the arrow denotes the flow when we *lower* $\log P$, following the standard convention.

We see that when $g_*^2 = -b_1/b_2$, the value of g^2 becomes independent of the scale $\log P$ at which we observe the system, at least to this order of approximation. (The value of the coupling at the fixed point is conventionally denoted by attaching $*$ as a subscript or a superscript.) Such a value is called a *renormalization group fixed point*. We are supposing that g^2 is an inherently positive quantity, so this requires that $g_*^2 = -b_1/b_2 > 0$. Furthermore, to be really sure about the existence of the fixed point, g_*^2 needs to be sufficiently small. We also have a much more trivial fixed point when $g^2 = 0$: this is called a *free fixed point* or a *Gaussian fixed point*.

Depending on whether the fixed point is approached in the ultraviolet (i.e. at high energy or at the small scale) or in the infrared (i.e. at low energy or at long distance), the fixed points are called ultraviolet or infrared, respectively. In Fig. 13, we also indicated the types of the fixed points in the infrared and in the ultraviolet, assuming that g_*^2 is very small. The entries marked by ? are for the cases where the renormalization group flow takes the value of g^2 out side of the region where the perturbation theory is valid.

When the coupling g_*^2 stays constant under the change of $\log P$, the system is scale invariant. Such a system is often invariant under a bigger symmetry called the conformal symmetry, and is called a conformal field theory.

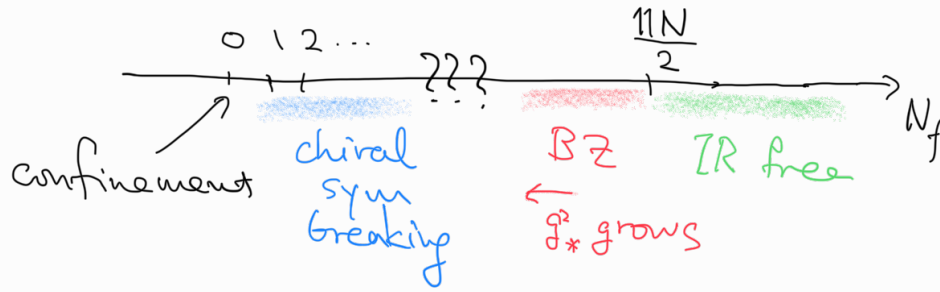


Figure 14: Schematic picture of the phase structure of the infrared limit of QCD.

5.4 The Banks-Zaks IR fixed point

Let us now study some example. The two-loop running of the $SU(N)$ gauge theory with N_f quarks in the fundamental representations was first computed by [Cas74, Jon74, BM74], and has the form:

$$\frac{\partial}{\partial \log P} g^2 = \frac{g^4}{8\pi^2} \left(-\frac{11}{3}N + \frac{2}{3}N_f \right) + \frac{g^6}{(8\pi^2)(16\pi^2)} \left(-\frac{34}{3}N^2 + \left(\frac{20}{3}N + 4\frac{N^2-1}{2N} \right) \frac{N_f}{2} \right) + O(g^8). \quad (5.30)$$

Exercise 5.8. Compute and confirm the coefficient of g^6 .

We take the limit $N, N_f \rightarrow \infty$ keeping $N_f/N = 11/2 - \epsilon$ with a very small ϵ fixed. We find that it has a weakly-coupled infrared fixed point at

$$\frac{g_*^2 N}{16\pi^2} = \frac{4}{75} \epsilon + O(\epsilon^2). \quad (5.31)$$

This was first noted in [BZ82] and is known as the Banks-Zaks fixed point after the discoverer.

As we will study in detail later, when $N_f = 0$, the $SU(N)$ gauge theory is believed to confine and develop a mass gap. When N_f is small, the QCD will break the chiral symmetry by having the nonzero vev $\langle \psi_\alpha \psi^\alpha \rangle \neq 0$. When $N_f \lesssim (11/2)N$, the long-distance limit is the Banks-Zaks fixed point, which is a weakly-coupled conformal field theory, and when $N_f > (11/2)N$, the system is infrared free with a logarithmically decreasing coupling. The range of N_f where the infrared limit is a nontrivial conformal field theory is known as *the conformal window*. As we lower N_f/N , the fixed point coupling g_*^2 grows, and eventually we lose perturbative control. This makes it hard to find the lower end of the conformal window.

Our world has $N = 3$, and depending on the energy scale, we usually think of N_f as 2 or 3. Experimentally, we know that we have the chiral symmetry breaking. So $N_f = 3$ is below the lower boundary. The upper boundary is at $N_f < (11/2)N$, i.e. $N_f \leq 16$. A lot of numerical efforts were put forward to determine the lower boundary. See [DeG15] for an extensive review from 2015. The author of that review says that the evidence that $N_f = 8$ is below the lower boundary is very strong, while $N_f = 12$ is definitely above the lower boundary.

5.5 The Litim-Sannino UV fixed point

A natural question then is whether there is a weakly-coupled ultraviolet fixed point. This question was only (relatively¹²) recently answered, in the affirmative by Litim and Sannino [LS14].¹³

The model is not too difficult. One starts from $SU(N)$ gauge theory with N_f quarks ψ_L^i and $\bar{\psi}_{Rj}$ in the fundamental representation, $i = 1, \dots, N_f$. We then introduce $N_f \times N_f$ complex scalar fields H_j^i . One just consider the Lagrangian

$$S = - \int d^4x \left[\frac{1}{2g^2} \text{tr} F_{\mu\nu} F^{\mu\nu} + \bar{\psi}_L D_\mu \sigma^\mu \psi_L + \psi_R D_\mu \sigma^\mu \bar{\psi}_R + \partial_\mu H^\dagger \partial_\mu H \right. \\ \left. + y \psi_L^i \psi_{R,j} H_i^j + c.c. + u \text{tr} (HH^\dagger)^2 + v (\text{tr} HH^\dagger)^2 \right] \quad (5.32)$$

The couplings are

$$\alpha_g := \frac{g^2 N}{16\pi^2}, \quad \alpha_y := \frac{y^2 N}{16\pi^2}, \quad \alpha_u := \frac{u N_F}{16\pi^2}, \quad \alpha_v := \frac{v N_F^2}{16\pi^2}. \quad (5.33)$$

As before, we define

$$N_F/N = 11/2 + \epsilon \quad (5.34)$$

and take the limit $N_F, N \rightarrow \infty$, keeping ϵ small and fixed.

The RG equations are, according to [LS14, LMS15],

$$\dot{\alpha}_g = \frac{4}{3} \epsilon \alpha_g + 25 \alpha_g^2 - \frac{121}{2} \alpha_g \alpha_y + O(\alpha^3), \quad (5.35)$$

$$\dot{\alpha}_y = \alpha_y (13 \alpha_y - 6 \alpha_g + O(\alpha^2)), \quad (5.36)$$

$$\dot{\alpha}_u = -11 \alpha_y^2 + 4 \alpha_u (\alpha_y + 2 \alpha_u) + O(\alpha^3), \quad (5.37)$$

$$\dot{\alpha}_v = 12 \alpha_u^2 + 4 \alpha_v (\alpha_v + 4 \alpha_u + \alpha_y) + O(\alpha^3). \quad (5.38)$$

Here the dot on the left hand side means the derivative with respect to the scale $\log P$.

The structure of the RG equations show that a consistent fixed point can be found at $\alpha^* \sim \epsilon$. From (5.36) one finds $\alpha_y^* = (6/13) \alpha_g^*$. Plugging this in to (5.35) one finds α_g^* and α_y^* . The equation (5.37) then determines α_u^* , which can be fed to (5.38) to fix α_v^* . The results are:

$$\alpha_g^* = \frac{26}{57} \epsilon, \quad \alpha_y^* = \frac{4}{19} \epsilon, \quad \alpha_u^* = \frac{\sqrt{23} - 1}{19} \epsilon, \quad \alpha_v^* = -\frac{1}{19} (2\sqrt{23} - \sqrt{20 + 6\sqrt{23}}) \epsilon, \quad (5.39)$$

and this fixed point is in the ultraviolet; the flow is visualized in Fig. 15.

¹²In the sense that the answer was not known when I was a graduate student around 2005. I remember being asked by Prof. Yanagida whether such a UV fixed point exists. I looked for various supersymmetric theories, because I never dealt with non-supersymmetric theories back then. I did not find any. Indeed, later, Intriligator and Sannino found that you cannot have weakly-coupled UV fixed points in supersymmetric theories [IS15].

¹³I might be writing too many personal recollections here, but there is another point I want to make. The authors of this paper has been pursuing the idea called the asymptotic safety of gravity, in which they hope for the existence of an ultraviolet fixed point of gravity. If it were realized, you do not need string theory or loop quantum gravity to deal with quantum gravity, and the ordinary QFT analysis would suffice. Personally I thought this approach was somewhat misguided, and I still think it is. That said, while trying to look for pieces of evidence, they analyzed ordinary gauge theories, and found this class of weakly-coupled UV fixed points. This gave me a cautionary tale: interesting results can come out even from a misguided motivation. It is often useful, therefore, to let other people pursue what they like, independent of whether you like it or not.

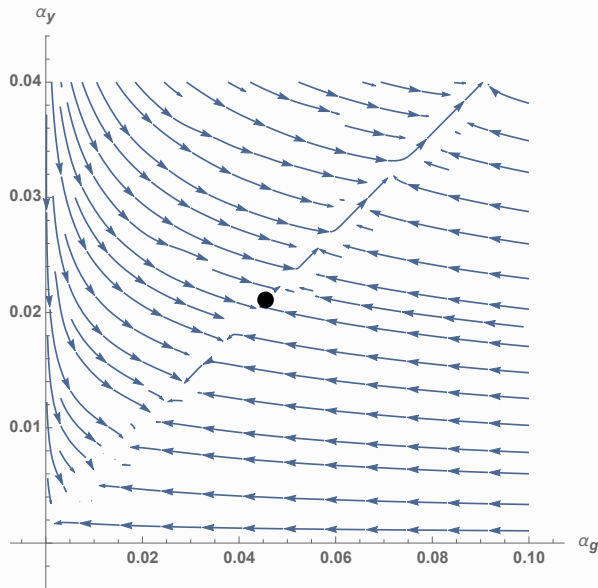


Figure 15: Renormalization group flow for α_g and α_y . We chose $\epsilon = 0.1$; the arrow shows the direction toward the infrared; the black blob is the fixed point.

That $\alpha_v < 0$ is somewhat worrisome, since this might make the potential for H unbounded below. To check that this does not happen, one rewrites

$$u \operatorname{tr} X^2 + v(\operatorname{tr} X)^2 = u \operatorname{tr}(X - (\operatorname{tr} X)/N_f)^2 + (v - u/N_f)(\operatorname{tr} X)^2 \quad (5.40)$$

where $X = HH^\dagger$. We can use (5.39) to check that $v^* - u^*/N_f = (\alpha_v^* - \alpha_u^*)N_f^2 > 0$. Therefore the scalar potential at the fixed point is bounded from below.

6 Qualitative discussions of strongly-coupled gauge theories

We saw above that $SU(N)$ gauge theory with N_f massless quarks in the fundamental representation is

- infrared free when $N_f > (11/2)N$,
- goes to the Banks-Zaks conformal fixed point when $N_f \lesssim (11/2)N$.

We also learned that the coupling g_* at the Banks-Zaks fixed point grows as we lower N_f .

What happens in the other extreme, when $N_f = 0$ (which is known as the pure gauge theory or the pure Yang-Mills theory) or when $N_f = 1, 2$ or 3 ? We do not have any perturbative control in the infrared limit. Here the experiment gives us important information, since our world has the strong force and the quarks, corresponding to $N = 3$ and $N_f = 2$ or 3 .

6.1 $N_f = 0$: color confinement and the mass gap

When $N_f = 0$, i.e. in the pure $SU(N)$ gauge theory, with the bare action

$$S = - \int \frac{1}{2g^2} \text{tr} F_{\mu\nu} F^{\mu\nu}, \quad (6.1)$$

the system is believed to *confine its color* and *generate the mass gap*. This means the following.

Color confinement: First, recall the situation in the more familiar quantum electrodynamics. In the classical theory of Maxwell fields and electrons, the electron field couple to the Maxwell fields. In other words, the electron fields are electrically charged. In the quantum theory of Maxwell fields and electrons, there is a single electron state in the Hilbert space, which has finite energy and has an electric charge.

Now, let us consider the classical pure $SU(N)$ theory. The gauge fields have self couplings, since $F_{\mu\nu} = \partial_\mu A_\nu - \partial_\nu A_\mu + i[A_\mu, A_\nu]$ and therefore $\text{tr} F_{\mu\nu} F^{\mu\nu}$ contains cubic and quartic couplings. What couples to $SU(N)$ gauge fields are said to have color charges. After a naive quantization, the gauge fields give the gluons. And indeed in the high energy limit, they are known to interact weakly. But it is believed that it costs infinite amount of energy to isolate a single colored particle. In other words, in the Hilbert space of the theory, there is no finite energy state which has the color charge. This is the color confinement.

Generation of the mass gap: In general, a one-particle excitation over a relativistic vacuum is characterized by a 4-momentum of the form $E^2 - p^2 = m^2$, where m is the rest mass of that excitation. In the quantum Maxwell theory, there definitely is the photon, which is massless. In the quantum $SU(N)$ gauge theory, the gluon has colors and cannot be isolated. The lowest-mass excitation above the vacuum is known as the glue ball, and it is believed to be massive, despite the fact that there is no mass term in the original Lagrangian above. This is the *generation of the mass gap*.

Here are two anecdotes:

1. One of the Clay Millennium problems is to mathematically construct the pure Yang-Mills theory and prove that there is the mass gap [JW06]. A successful proof will give you a million dollars.
2. When Yang and Mills originally came up with the Yang-Mills theory, Yang gave a seminar at IAS. He wrote his experience in the comment section of his collected works, see pp. 19–20 of [Yan05]. Let me quote the most interesting part:

Oppenheimer invited me to return to Princeton for a few days in late February to give a seminar on our work. Pauli was spending the year in Princeton, and he was deeply interested in symmetries and interactions ... Soon after my seminar began, when I had written down on the blackboard,

$$(\partial_\mu - ieB_\mu)\psi$$

Pauli asked, “What is the mass of this field B ?” I said we did not know. Then I resumed my presentation, but soon Pauli asked the same question again. I said something to the effect that that was a very complicated problem, we had worked on it and had come to no definite conclusions. I still remember his repartee: “That is not sufficient excuse.” I was so taken aback that I decided, after a few moments’ hesitation, to sit down. There was general embarrassment. Finally Oppenheimer said, “We should let Frank proceed.” I then resumed, and Pauli did not ask any more questions during the seminar.

Exercise 6.1. Prove that pure Yang-Mills theory has a mass gap.

6.2 $N_f = 1, 2, \dots$: chiral symmetry breaking and the U(1) problem

Before proceeding, we note that two phenomena, the color confinement and the generation of the mass gap, are logically independent. Indeed, with $N_f = 2, 3$, it is believed that the system confines without the mass gap. To understand this we need to know the concept of the chiral symmetry.

Let us recall the classical action:

$$S = - \int \left[\frac{1}{2g^2} \text{tr} F_{\mu\nu} F^{\mu\nu} + \bar{\psi}_{ai\dot{\alpha}} \sigma_{\dot{\alpha}\alpha}^{\mu} D_{\mu} \psi_{\alpha}^{ai} + \tilde{\psi}_{\dot{\alpha}}^{au} \sigma_{\dot{\alpha}\alpha}^{\mu} D_{\mu} \tilde{\psi}_{au\alpha} \right] \quad (6.2)$$

where we put indices explicitly to fermions. First $\alpha, \dot{\alpha}$ are spinor indices. Second, $a = 1, \dots, N$ are indices for the color (=gauge) indices. This means that ψ is in the fundamental representation (a column vector of N elements, say), and $\tilde{\psi}$ is in the antifundamental representation (a row vector of N elements). $\bar{\psi}$ is then in the antifundamental and $\tilde{\bar{\psi}}$ is in the fundamental.

Third, $i = 1, \dots, N_f$ and $u = 1, \dots, \tilde{N}_f$ are for the distinct flavors of quarks. For the moment we consider introducing different numbers of quarks for ψ and $\tilde{\psi}$. We can compute the gauge anomaly as explained in Sec. 3.4. We consider a U(1) subgroup in the SU(N) group corresponding to the infinitesimal generator

$$X = \text{diag}(q_1, \dots, q_N), \quad \sum q_a = 0. \quad (6.3)$$

Under this U(1) subgroup, the a -th component of ψ has charge q_a and the a -th component of $\tilde{\psi}$ has charge $-q_a$. Then the total gauge anomaly is

$$N_f \sum (q_a)^3 + \tilde{N}_f \sum (-q_a)^3, \quad (6.4)$$

which is zero if and only if $N_f = \tilde{N}_f$. We want the gauge anomaly to be absent, so we assume $N_f = \tilde{N}_f$ in the following. In this case, we can stop distinguishing two types of indices i and u and we can combine two Weyl spinors ψ_{α}^{ai} and $\tilde{\psi}_{\dot{\alpha}}^{ai}$ into a Dirac spinor $\Psi^{ai} = (\psi_{\alpha}^{ai}, \tilde{\psi}_{\dot{\alpha}}^{ai})$.

We now consider the field redefinition of the form

$$\psi_{\text{old}}^i \rightarrow \psi_{\text{new}}^i := \mathcal{U}_j^i \psi_{\text{old}}^j, \quad \tilde{\psi}_{\text{old}}^u \rightarrow \tilde{\psi}_{\text{new}}^u := \psi_{\text{old}}^v \mathcal{V}^{-1v}_u, \quad (6.5)$$

where $\mathcal{U} \in U(N_f)$ and $\mathcal{V} \in U(N_f)$ are spacetime-independent matrices. This is clearly a symmetry of the classical Lagrangian above.

Historically, people originally used the Dirac spinor Ψ^{ai} , with which the transformation

$$\Psi_{\text{old}}^i \rightarrow \Psi_{\text{new}}^i := \mathcal{W}_j^i \Psi_{\text{old}}^j, \quad (6.6)$$

looks more obvious. This corresponds to the special case $(\mathcal{U}, \mathcal{V}) = (\mathcal{W}, \mathcal{W})$. \mathcal{W} are said to form the diagonal subgroup of the group of $(\mathcal{U}, \mathcal{V})$.

The transformations

$$(\mathcal{U}, \mathcal{V}) \in U(N_f) \times U(N_f) \quad (6.7)$$

is known as the chiral flavor symmetry, and the subgroup

$$\mathcal{W} \in U(N_f)_{\text{diagonal}} \subset U(N_f) \times U(N_f) \quad (6.8)$$

is called the the flavor symmetry. We will soon see that the non-diagonal part of the $U(1)$ symmetry is in fact anomalous and is not the symmetry of the theory, but for the moment we ignore this important subtlety.

Chiral symmetry breaking: For a sufficiently small N_f , this system is believed to exhibit *color confinement* and *chiral symmetry breaking*. We already introduced the color confinement, so let us concentrate on the chiral symmetry breaking. The essence is that the strong interaction in the infrared is believed to produce a non-zero eigenvalue for the fermion bilinear

$$\langle \psi_\alpha^{ai} \tilde{\psi}_{au}^\alpha \rangle \propto \Lambda_{\text{dyn}}^3 \delta_u^i. \quad (6.9)$$

This is also known as the *chiral condensate*, and can roughly be thought of as the analogue of the Cooper pair in the superconductor.

Note that the chiral symmetry acts on this condensate as

$$\langle \psi_{\text{old}}^i \tilde{\psi}_u^{\text{old}} \rangle \mapsto \langle \psi_{\text{new}}^i \tilde{\psi}_u^{\text{new}} \rangle = \mathcal{U}_j^i \langle \psi_{\text{old}}^j \tilde{\psi}_v^{\text{old}} \rangle \mathcal{V}^{-1v}_u, \quad (6.10)$$

and therefore is invariant only under $(\mathcal{U}, \mathcal{V}) = (\mathcal{W}, \mathcal{W})$, because of the presence of δ_u^i . In other words, the nonzero condensate breaks the chiral symmetry $U(N_f) \times U(N_f)$ to the non-chiral subgroup $U(N_f)_{\text{diag}}$.

The $U(1)$ problem: If this were the whole story, we would obtain $2N_f^2 - N_f^2 = N_f^2$ massless Nambu-Goldstone bosons associated to the symmetry breaking, because we have $2N_f^2$ generators for $U(N_f) \times U(N_f)$ which is broken to $U(N_f)$ with N_f^2 generators. However, in nature where $N_f = 2$, we observed only $N_f^2 - 1 = 3$ massless Nambu-Goldstone bosons, which were identified as the pions π^+ , π^0 and π^- , which corresponded to broken generators of $SU(N_f) = SU(2)$. (Note that they are actually massive, about $\sim 135\text{MeV}$, but this can be accounted for from the correction coming from the fact that the quarks are not massless. The problem is that the next lightest particle with the correct quantum number, η , is too heavy.) What happened to the would-be Nambu-Goldstone boson for the $U(1)$ part of the broken generator?

Let us be more explicit about the structure of the symmetry. We now separate $U(N_f) \times U(N_f)$ as $(U(1) \times SU(N_f)) \times (U(1) \times SU(N_f))$. The combination of $U(1) \times U(1)$, which assigns charge $+1$ to ψ and -1 to $\tilde{\psi}$, is called the *baryonic symmetry*, and is an actual symmetry. This combination is called the baryonic symmetry since it assigns zero charge to a meson operator $\psi\tilde{\psi}$ whereas it assigns charge $+N$ to the baryon operator

$$\epsilon_{a_1 a_2 \dots a_N} \psi_{\alpha_1}^{a_1} \psi_{\alpha_2}^{a_2} \dots \psi_{\alpha_N}^{a_N}. \quad (6.11)$$

The other combination of $U(1) \times U(1)$, which assigns charge $+1$ to both ψ and $\tilde{\psi}$, is called the axial $U(1)$ ‘symmetry’, but it is simply not a symmetry of the system in the quantum theory. We will detail why this is the case in the next subsection, but let us accept it for a while.

Such a ‘symmetry’ is called anomalous.¹⁴ The Nambu-Goldstone theorem does not apply since an anomalous ‘symmetry’ is not really a symmetry, and therefore we do not have to worry about the missing Nambu-Goldstone boson for the axial $U(1)$ symmetry. This is the famous solution to the $U(1)$ problem by ’t Hooft [tH76, tH86].¹⁵ Summarizing, we have the situation:

$$\begin{array}{ll} \text{symmetry of the classical Lagrangian:} & U(N_f) \times U(N_f) \\ & \downarrow \\ \text{symmetry of the quantum theory:} & U(1)_{\text{baryon}} \times SU(N_f) \times SU(N_f) \\ & \downarrow \\ \text{symmetry of the infrared limit:} & U(1)_{\text{baryon}} \times SU(N_f)_{\text{diagonal}}. \end{array} \quad (6.12)$$

where $N_f^2 - 1$ generators are spontaneously broken by the chiral condensate in the second arrow, producing $N_f^2 - 1$ Nambu-Goldstone bosons in the process.

6.3 The $U(1)$ problem

We said above that the chiral $U(1)$ ‘symmetry’ is not a symmetry to start with. Let us understand why this is the case.

6.3.1 First derivation

To see this, we repeat the argument in Sec. 3.4. We pick an arbitrary $U(1)$ subgroup of $U(N_f) \times U(N_f)$ corresponding to the generator

$$\text{diag}(q_1, q_2, \dots, q_{N_f}), \quad \text{diag}(\tilde{q}_1, \tilde{q}_2, \dots, \tilde{q}_{N_f}) \quad (6.13)$$

for $(\mathcal{U}, \mathcal{V}) \in U(N_f) \times U(N_f)$, respectively, and would like to ask which $U(1)$ subgroup is actually a symmetry.

¹⁴This is a confusing terminology, since an anomalous ‘symmetry’ is not really a symmetry. But it is hard to change the established terminology.

¹⁵According to <http://www.staff.science.uu.nl/~hooft101/ap.html>, ’t is pronounced as [ət]. Therefore with a indefinite article we should use *an ’t Hooft loop* not *a ’t Hooft loop*, and similarly *the ’t Hooft anomaly matching condition* should be pronounced with [ð̥i] instead of [ð̥ə].

Let us introduce a magnetic flux B_z for this $U(1)$ subgroup with $\int B_z dx dy = 2\pi$. More explicitly, we introduce a magnetic flux for $U(N_f) \times U(N_f)$ of the form

$$\text{diag}(q_1, q_2, \dots, q_{N_f}) B_z \quad \text{for } U(N_f) \text{ acting on } \psi \quad (6.14)$$

and

$$\text{diag}(\tilde{q}_1, \tilde{q}_2, \dots, \tilde{q}_{N_f}) B_z \quad \text{for } U(N_f) \text{ acting on } \bar{\psi}. \quad (6.15)$$

This leads to $(q_1 + \dots + q_{N_f})$ massless *left-moving* 2d fermions in the $z - t$ plane charged under gauge $SU(N)$ from ψ_α , and similarly to $(\tilde{q}_1 + \dots + \tilde{q}_{N_f})$ massless *right-moving* 2d fermions in the $z - t$ plane charged under gauge $SU(N)$ from $\bar{\psi}_{\dot{\alpha}}$. This has gauge anomaly and therefore inconsistent unless there are equal number of left-moving and right-moving fermions, i.e.

$$q_1 + \dots + q_{N_f} = \tilde{q}_1 + \dots + \tilde{q}_{N_f}. \quad (6.16)$$

The $SU(N_f) \times SU(N_f)$ part automatically satisfy the constraint above, since we have

$$q_1 + \dots + q_{N_f} = 0 = \tilde{q}_1 + \dots + \tilde{q}_{N_f} \quad (6.17)$$

separately. We now set $q_1 = \dots = q_{N_f} =: q$ and $\tilde{q}_1 = \dots = \tilde{q}_{N_f} =: \tilde{q}$ to isolate the $U(1) \times U(1)$ part. We see that there is a gauge anomaly unless $q = \tilde{q}$. This means that the diagonal $U(1) \subset U(1) \times U(1)$ has no problem, whereas a magnetic field in the anti-diagonal $U(1)$ where $q = -\tilde{q}$ causes a gauge anomaly.

6.3.2 Second derivation

In the above derivation, we saw the appearance of the problematic gauge anomaly of $SU(N)$ when one introduces a nonzero background gauge field for the chiral $U(1)$ symmetry. Let us see how the problem manifests itself when we consider a nonzero gauge field for $SU(N)$.

Let us first consider a single Weyl fermion ψ of charge $+1$ under $U(1)$. Let us introduce the background $U(1)$ gauge field on Euclidean \mathbb{R}^4 parameterized by x, y, z, w , such that only the components F_{xy} and F_{zw} are nonzero, with

$$\int dx dy \frac{F_{xy}}{2\pi} = n_{xy}, \quad \int dz dw \frac{F_{zw}}{2\pi} = n_{zw}. \quad (6.18)$$

We first regard x, y directions as internal; we saw there are n_{xy} zero modes of \not{D} for the x, y directions. Then we effectively have n_{xy} left-moving complex 2d fermions on the zw plane. We apply the same argument again on the z, w directions; there are n_{zw} zero modes of \not{D} for the z, w directions, for each left-moving complex 2d fermion. In total, we have $n_{xy} n_{zw}$ zero modes of the 4d Dirac operator \not{D} .

We can write this combination $n_{xy} n_{zw}$ in a more covariant way:

$$n_{xy} n_{zw} = \int d^4x \frac{1}{8} \frac{1}{(2\pi)^2} \epsilon_{xyzw} F_{\mu\nu} F_{\mu\nu}, \quad (6.19)$$

and in this way it is known to generalize:

$$((\text{effective}) \# \text{ of the zero modes of } \not{D}) = \int d^4x \frac{1}{8} \frac{1}{(2\pi)^2} \epsilon_{\mu\nu\rho\sigma} F_{\mu\nu} F_{\rho\sigma}, \quad (6.20)$$

for arbitrary configuration of F , not just for configurations of the particular form (6.18). This is the index theorem.

What does a 4d zero mode do? We have a normalizable solution to

$$\sigma^\mu D_\mu \psi = 0 \quad (6.21)$$

in other words

$$(\partial_t \psi(t, x, y, z) + \sum_{i=1,2,3} \sigma^i D_i) \psi(t, x, y, z) = 0. \quad (6.22)$$

We can expand the x, y, z directions using the eigenmodes of $\sigma^i D_i$, which is now time dependent. Assuming that the change is very slow, we see

$$(\partial_t \psi_n(t) + E_n(t)) \psi_n(t). \quad (6.23)$$

Integrating it from $t = 0$ to $t \rightarrow \pm\infty$, we see that

$$\psi_n(t) \sim e^{-\int_0^t E_n(t) dt} \quad (6.24)$$

which converges only if $E_n(t) > 0$ for $t \gg 0$ and $E_n(t) < 0$ for $t \ll 0$. Therefore the zero modes correspond to the energy eigenvalues crossing from negative to positive. This brings up a fermionic mode from the Dirac sea. Since each new fermionic mode carries the U(1) charge, we have

$$Q(t = +\infty) - Q(t = -\infty) = ((\text{effective}) \# \text{ of the zero modes of } \not{D}). \quad (6.25)$$

Combining with (6.20), we have

$$Q(t = +\infty) - Q(t = -\infty) = \int d^4x \frac{1}{8} \frac{1}{(2\pi)^2} \epsilon_{\mu\nu\rho\sigma} F_{\mu\nu} F_{\rho\sigma}. \quad (6.26)$$

In fact there is a local version of this equation:

$$\partial^\mu j_\mu = \frac{1}{8} \frac{1}{(2\pi)^2} \epsilon_{\mu\nu\rho\sigma} F_{\mu\nu} F_{\rho\sigma} \quad (6.27)$$

where j_μ is the U(1) current. This form of the equation can be derived directly by a perturbation theory one-loop computation; the coefficient is easier to fix in the approach taken in this section in my opinion. Technicality aside, all this means that the U(1) current is not conserved and the U(1) is not really a symmetry.

We can now repeat the same discussion in the case of SU(N) gauge theory with N_f flavors. First consider the case $N_f = 1$. We have ψ_α in the fundamental of SU(N) and $\tilde{\psi}_\alpha$ in the antifundamental of SU(N). We assign U(1) charge q for ψ_α and $-\tilde{q}$ for $\tilde{\psi}_\alpha$. The same argument as above shows that the U(1) current is anomalous in general:

$$\partial_\mu j^\mu = (q - \tilde{q}) \frac{1}{8} \frac{1}{(2\pi)^2} \epsilon_{\mu\nu\rho\sigma} \text{tr} F_{\mu\nu} F_{\rho\sigma}. \quad (6.28)$$

Then the diagonal U(1) symmetry for which $q = \tilde{q}$ is actually a symmetry while the chiral U(1) symmetry for which $q = -\tilde{q}$ is *not* a symmetry.

For the more general case of $N_f > 1$, we again pick the U(1) subgroup specified by q_1, \dots, q_{N_f} and $\tilde{q}_1, \dots, \tilde{q}_{N_f}$. Then the anomaly equation is

$$\partial_\mu j^\mu = \left(\sum_i q_i - \sum_i \tilde{q}_i \right) \frac{1}{8} \frac{1}{(2\pi)^2} \epsilon_{\mu\nu\rho\sigma} \text{tr} F_{\mu\nu} F_{\rho\sigma}. \quad (6.29)$$

Therefore we come to the same conclusion as in the first method: $U(1)_{\text{baryon}} \times SU(N_f) \times SU(N_f)$ is actually a symmetry, while $U(1)_{\text{chiral}}$ is not.

6.3.3 Instanton number

In the equation above, the integral of the left hand side is $Q(t = +\infty) - Q(t = -\infty)$ and is an integer. Then the right hand side should also be an integer. This means that, for an $SU(N)$ gauge field,

$$I := \int d^4x \frac{1}{8} \frac{1}{(2\pi)^2} \epsilon_{\mu\nu\rho\sigma} \text{tr} F_{\mu\nu} F_{\rho\sigma} \quad (6.30)$$

should be an integer. This integer is known to be the instanton number. Is there a way to understand that it is an integer?

Let us recall a simpler question, which is the Dirac quantization of the monopole number. This was

$$\int d^2x \frac{1}{2} \frac{1}{2\pi} \epsilon_{\mu\nu} F_{\mu\nu} \quad (6.31)$$

for a U(1) gauge field, and the integral is over a sphere. This was shown as follows. We separate the integral to the sum of the integrals over the northern and the southern hemisphere. On each hemisphere, we use

$$\frac{1}{2} \epsilon_{\mu\nu} F_{\mu\nu} = \epsilon_{\mu\nu} \partial_\mu A_\nu \quad (6.32)$$

to partially integrate. Then

$$\int d^2x \frac{1}{2} \frac{1}{2\pi} \epsilon_{\mu\nu} F_{\mu\nu} = \int_{\text{equator}} d\theta \frac{1}{2\pi} (A_\theta^{\text{north}} - A_\theta^{\text{south}}) = \int_{\text{equator}} d\theta \frac{1}{2\pi} \partial\chi = \chi(\theta = 2\pi) - \chi(\theta = 0) \quad (6.33)$$

where $e^{i\chi(\theta)}$ is the gauge transformation on the equator.

The approach to the case above is similar. We separate \mathbb{R}^4 into the region $t < 0$ and $t > 0$. We then use the equation

$$\frac{1}{8} \frac{1}{(2\pi)^2} \epsilon_{\mu\nu\rho\sigma} \text{tr} F_{\mu\nu} F_{\rho\sigma} = \frac{1}{2} \frac{1}{(2\pi)^2} \epsilon_{\mu\nu\rho\sigma} \partial_\mu \text{tr} (A_\nu \partial_\rho A_\sigma + \frac{2}{3} A_\nu A_\rho A_\sigma) \quad (6.34)$$

to partially integrate, so that

$$I = \frac{1}{2} \frac{1}{(2\pi)^2} \int_{|x|=R} \left[d^3x \epsilon_{\nu\rho\sigma} \text{tr} (A_\nu \partial_\rho A_\sigma + \frac{2}{3} A_\nu A_\rho A_\sigma) - \text{tr} (A'_\nu \partial_\rho A'_\sigma + \frac{2}{3} A'_\nu A'_\rho A'_\sigma) \right] \quad (6.35)$$

where A_ν and A'_ν are the gauge fields on $t < 0$ and $t > 0$, respectively.

Exercise 6.2. Check this equation (and correct my typos).

Now they are related by the gauge transformation

$$A'_\nu = gA_\nu g^{-1} + g\partial_\nu g^{-1}. \quad (6.36)$$

After plugging this in to the equation above, and after a somewhat lengthy computation when the dust settles, one finds

$$I = \frac{1}{24\pi^2} \int_{t=0} d^3x \epsilon_{\nu\rho\sigma} \text{tr}(g\partial_\nu g^{-1})(g\partial_\rho g^{-1})(g\partial_\sigma g^{-1}) \quad (6.37)$$

When $G = \text{SU}(2)$, g is a map from $\mathbb{R}^3 = \{t = 0\}$ to $S^3 \simeq \text{SU}(2)$. We can add a point at infinity to think that \mathbb{R}^3 is essentially an S^3 , so g is a map from S^3 to S^3 . I computes the winding number, which is an integer. When $G = \text{SU}(N)$, I computes a generalized version of this winding number and still is an integer.

Exercise 6.3. Compute the winding number explicitly for the identity map $S^3 \rightarrow S^3$ using this formula, and check that $I = 1$.

We now have another way to see the chiral anomaly. A $\text{SU}(N)$ background with nonzero instanton number I requires a global gauge transformation on some constant-time slice $g : \mathbb{R}^3 \rightarrow \text{SU}(N)$ with the winding number I . This creates the chiral charge $(\sum_i q_i - \sum_i \tilde{q}_i)I$. A gauge transformation is a redundancy in the description, and should not change a genuine physical observable. Therefore, when $\sum_i q_i - \sum_i \tilde{q}_i \neq 0$, the $\text{U}(1)$ subgroup is not a genuine symmetry.

6.4 Chiral Lagrangian

Assuming that there is a chiral condensate $\langle \psi_\alpha^{ai} \tilde{\psi}_{au}^\alpha \rangle \neq 0$, we can proceed to the study of the Goldstone bosons. We saw above that the actual symmetry is $\text{U}(1)_{\text{baryon}} \times \text{SU}(N_f) \times \text{SU}(N_f)$, which is then broken to $\text{U}(1)_{\text{baryon}} \times \text{SU}(N_f)$ by the vev

$$\langle \psi_\alpha^{ai} \tilde{\psi}_{au}^\alpha \rangle \propto \Lambda^3 \delta_u^i \quad (6.38)$$

There should then be $N_f^2 - 1$ broken symmetry directions, which should be massless. How do we describe its dynamics?

We can start by a simpler example. Consider a theory of a single complex scalar field ϕ with a $\text{U}(1)$ symmetry, with the potential $V(\phi) = \lambda(|\phi|^2 - v^2)^2$. The lowest-energy point is $\langle \phi \rangle = v e^{i\theta}$ and the fluctuation along θ is the Nambu-Goldstone mode associated to the breaking of the $\text{U}(1)$ symmetry. We note that all possible vev's are obtained by starting from one particular value $\langle \phi \rangle = v$ and applying the $\text{U}(1)$ transformation $\phi \mapsto e^{i\theta} \phi$.

The action of the Nambu-Goldstone mode is obtained by simply restricting the original Lagrangian

$$S = \int d^4x \partial_\mu \bar{\phi} \partial^\mu \phi + V(\phi) \quad (6.39)$$

to the those ϕ of the form $\phi = v e^{i\theta}$, resulting in

$$S = \int d^4x \partial_\mu \bar{\phi} \partial^\mu \phi = \int d^4x v^2 \partial_\mu \theta \partial^\mu \theta, \quad (6.40)$$

which describes a massless scalar field.

Let us now come back to the case of the chiral symmetry breaking of QCD. The Goldstone mode corresponds to the fluctuation of the chiral condensate (6.38)

$$\langle \psi_\alpha^{ai} \tilde{\psi}_{au}^\alpha \rangle \propto \Lambda^3 U_u^i \quad (6.41)$$

where U_u^i is obtained by δ_u^i by the action of the symmetry (6.10). We see that U_u^i can be an arbitrary unitary matrix of determinant 1.

We would like to write an action S of U_u^i . This needs to be symmetric under $SU(N_f) \times SU(N_f)$ acting on the indices i and u . The simplest term one can write has two derivatives, and the next has four derivatives:

$$S = \int d^4x \left[\frac{F_\pi^2}{4} \text{tr}(\partial_\mu U \partial^\mu U^\dagger) + a \text{tr}(\partial_\mu U \partial^\mu U^\dagger)^2 + \dots \right] \quad (6.42)$$

Here F_μ has dimension 1 and a is dimensionless. Note that one cannot write a potential term for U , since any monomial of U and U^\dagger symmetric under $SU(N_f) \times SU(N_f)$ is in fact a constant. This is as it should be, since the Goldstone modes should be massless. This Lagrangian is known as the chiral Lagrangian.

Let us be more specific and consider $N_f = 2$. It is customary to parameterize

$$U = \exp\left(\frac{i}{F_\pi} \begin{pmatrix} \pi^0(x) & \sqrt{2}\pi^+(x) \\ \sqrt{2}\pi^-(x) & -\pi^0(x) \end{pmatrix}\right), \quad (6.43)$$

where π^0 and π^\pm are pions. Then we have

$$\begin{aligned} \frac{F_\pi^2}{2} \text{tr}(\partial_\mu U \partial^\mu U^\dagger) &= \frac{1}{2} \partial_\mu \pi^0 \partial^\mu \pi^0 + \partial_\mu \pi^+ \partial^\mu \pi^- \\ &+ \frac{1}{6F_\pi^2} ((\pi^0 \partial_\mu \pi^0 + \pi^+ \partial_\mu \pi^- + \pi^- \partial_\mu \pi^+)^2 - ((\pi^0)^2 + 2\pi^+ \pi^-) (\partial_\mu \pi^0 \partial^\mu \pi^0 + 2\partial_\mu \pi^+ \partial^\mu \pi^-)) \\ &+ \dots \end{aligned} \quad (6.44)$$

Exercise 6.4. Perform this expansion.

Therefore the single $SU(N_f) \times SU(N_f)$ -invariant term $\text{tr } \partial_\mu U \partial^\mu U^\dagger$ contains the standard kinetic terms of the pions, together with a lot of interaction terms. The value F_π cannot be computed from first principles at this rough level of understanding, but we already know that the relative size of various four-point couplings among pions, up to a single constant F_π . We also know that the four-point interactions are all of order P^2/F_π^2 where P is the typical momentum scale.

What would be a contribution to the four-point couplings to the more complicated term in the chiral Lagrangian in (6.42), e.g. the term proportional to a ? A gedanken computation shows that any such terms would have four derivatives, and therefore carries a factor P^4/F_π^4 . This means that when $P \lll F_\pi$, the interaction processes are dominated by the contributions of the order P^2/F_π^2 we saw above. We can summarize our finding by saying that the behavior pion interactions in the low-momentum limit $P \lll F_\pi$ is completely determined by the single term $(F_\pi^2/4) \text{tr } \partial_\mu U \partial^\mu U$ alone.

So far we only considered the case when the quarks are massless. But we can use this formalism to study what happens when quarks are massive. Again let us say $N_f = 2$. We can introduce the mass term

$$m_1 \psi_\alpha^{a,i=1} \psi_{a,u=1}^\alpha m_2 \psi_\alpha^{a,i=2} \psi_{a,u=2}^\alpha + c.c. \quad (6.45)$$

to the original QCD Lagrangian (where we could call $m_1 = m_u$ and $m_2 = m_d$ are the masses of the up and down quarks). Using (6.41) we rewrite it as a term

$$m_1 \Lambda^3 U_{u=1}^{i=1} + m_2 \Lambda^3 U_{u=2}^{i=2} + c.c. = \Lambda^3 \text{tr} \left(\begin{pmatrix} m_1 & \\ & m_2 \end{pmatrix} U \right) + c.c. \quad (6.46)$$

Expanding in terms of π^0 and π^\pm , we find that it leads to the Lagrangian of the form

$$\frac{1}{2} \partial_\mu \pi^0 \partial^\mu \pi^0 + \frac{(m_1 + m_2) \Lambda^3}{F_\pi^2} ((\pi^0)^2 + 2\pi^+ \pi^-) + \dots \quad (6.47)$$

which means that the masses of π^0 and π^\pm are the same at this order and furthermore

$$m_\pi^2 \propto \frac{\Lambda^3 (m_1 + m_2)}{F_\pi^2}, \quad (6.48)$$

i.e. the pion masses are proportional to the square root of the quark mass.

6.5 Baryon as a soliton in the chiral Lagrangian

Our chiral Lagrangian only contains massless pions, bound states of the form $\psi \tilde{\psi}$. How do we describe a baryon, which is of the form

$$B := \psi^{a_1} \psi^{a_2} \dots \psi^{a_N} \epsilon_{a_1 \dots a_N} ? \quad (6.49)$$

Note that it has charge N under the $U(1)$ baryon symmetry.

To get a clue, consider $U(1)_{\text{baryon}}$ and $SU(N_f)_L \times SU(N_f)_R$. Repeating the computation of Sec. 6.3 but with the role of N_f and N exchanged, we can show that

$$\partial_\mu j_{\text{baryon}}^\mu = N \frac{1}{8} \frac{1}{(2\pi)^2} \epsilon_{\mu\nu\rho\sigma} (\text{tr } F_{\mu\nu}^L F_{\rho\sigma}^L - \text{tr } F_{\mu\nu}^R F_{\rho\sigma}^R) \quad (6.50)$$

where F^L and F^R are background gauge fields for $SU(N_f)_L \times SU(N_f)_R$ symmetry. This means that a gauge transformation in $SU(N_f)_L$ given by $g : \mathbb{R}^3 \rightarrow SU(N_f)$ of winding number w creates the baryon charge Nw , whereas the same in $SU(N_f)_R$ destroys the baryon charge Nw .

We now start from a pion configuration

$$\langle U_u^i \rangle = \delta_u^i, \quad (6.51)$$

on \mathbb{R}^3 at $t = 0$. This is the vacuum configuration and has baryon number 0, obviously. We can now apply an $SU(N_f)_L$ transformation $g : \mathbb{R}^3 \rightarrow SU(N_f)$ of winding number w to get the configuration

$$\langle U_u^i(x, y, z) \rangle = g_u^i(x, y, z). \quad (6.52)$$

We now see that this configuration has baryon number Nw .

In general, $\langle U_u^i(x, y, z) \rangle$ itself determines a map $U : \mathbb{R}^3 \rightarrow SU(N_f)$ and has a winding number. Therefore, any theory which has such a non-linear field U has the winding number as a topologically conserved quantum number. In the chiral Lagrangian of the QCD, we found that the winding number times N can be identified with the baryon number. This is in accord with the fact that every object in QCD with nonzero baryon number has baryon number which is a multiple of N , due to the fact that $\epsilon_{a_1 \dots a_N}$ has N indices.

This topological soliton is known as the Skyrmion, after Tony Skyrme who first considered such a configuration. The lowest energy configuration for the winding number 1 skyrmion needs to be found by solving the chiral Lagrangian (6.42). In the limit when the chiral Lagrangian solely consists of $\text{tr} \partial_\mu U \partial^\mu U^\dagger$, the Skyrmion tends to shrink, since for a configuration $U(x) := U(cx)$ one finds

$$\int d^3x \text{tr} \partial_i \tilde{U} \partial_i \tilde{U}^\dagger = \int d^3x c^2 \text{tr} (\partial_i U)(cx) (\partial_i U)(cx)^\dagger = c^{-1} \int d^3x \text{tr} \partial_i U \partial_i U^\dagger \quad (6.53)$$

which means that by shrinking it we can lower the energy arbitrarily.

In contrast, suppose there is a term containing four derivatives, such as $a \text{tr} (\partial_\mu U \partial^\mu U^\dagger)^2$. This term scales with the prefactor c^{+1} instead. Therefore there is a nonzero value of c where the combined energy

$$c^{-1} \int d^3x \frac{F_\pi^2}{4} \text{tr} \partial_i \tilde{U} \partial_i \tilde{U}^\dagger + c \int d^3x a \text{tr} (\partial_\mu U \partial^\mu U^\dagger)^2 \quad (6.54)$$

is minimized. In this way the size of the Skyrmion is fixed.

6.6 Wess-Zumino-Witten term

The chiral Lagrangian (6.42) is known to be quite useful in describing the pion dynamics, and even contains baryons as discussed above. But there are still a few points which are unsatisfactory.

One point is that the Lagrangian (6.42) is parity symmetric (in the sense that it contains no epsilon symbol), whereas the original QCD Lagrangian is not (in the sense that it uses Weyl spinors whose definition requires γ^5 which is essentially the epsilon symbol for spinors). We expect a term containing epsilon symbols in the chiral Lagrangian.

Another point is that the Lagrangian (6.42) does not explain the statistics of the baryons. Namely, a baryon contains N quark fields. Therefore, a baryon is a fermion if N is odd. How can a Skyrmion, which is just a nontrivial configuration of the bosonic field U , be a fermion when N is odd?

Luckily it is known that both problems can be solved by adding the so-called Wess-Zumino-Witten term [Wit83a, Wit83b]. Let us have a brief look at this term.

Let us consider a 2d version first. We already learned the concept of the winding number for $SU(N)$, given in (6.37), which we reproduce here:

$$I[g] = \frac{1}{24\pi^2} \int dt dx dy \epsilon_{\nu\rho\sigma} \text{tr}(g\partial_\nu g^{-1})(g\partial_\rho g^{-1})(g\partial_\sigma g^{-1}) \quad (6.55)$$

As discussed there, this is an integer, if it is integrated over the entire 3d space.

We now consider the following expression:

$$\Gamma_{\text{WZW}}[g] = \frac{1}{24\pi^2} \int_{y \geq 0} dt dx dy \epsilon_{\nu\rho\sigma} \text{tr}(g\partial_\nu g^{-1})(g\partial_\rho g^{-1})(g\partial_\sigma g^{-1}). \quad (6.56)$$

This is not an integer anymore. In addition, it has an interesting feature that the value modulo 1 only depends on the values of g at $y = 0$ and not on the values at $y > 0$.

To see this, let us say g_1 and g_2 are two maps $\{y \geq 0\} \rightarrow SU(N)$ such that their values at $y = 0$ agree. Then, we have

$$\begin{aligned} \Gamma_{\text{WZW}}[g_1] - \Gamma_{\text{WZW}}[g_2] &= \frac{1}{24\pi^2} \int_{y \geq 0} dt dx dy \epsilon_{\nu\rho\sigma} \text{tr}(g_1\partial_\nu g_1^{-1})(g_1\partial_\rho g_1^{-1})(g_1\partial_\sigma g_1^{-1}) \\ &\quad - \frac{1}{24\pi^2} \int_{y \geq 0} dt dx dy \epsilon_{\nu\rho\sigma} \text{tr}(g_2\partial_\nu g_2^{-1})(g_2\partial_\rho g_2^{-1})(g_2\partial_\sigma g_2^{-1}) \\ &= \frac{1}{24\pi^2} \int dt dx dy \epsilon_{\nu\rho\sigma} \text{tr}(g\partial_\nu g^{-1})(g\partial_\rho g^{-1})(g\partial_\sigma g^{-1}) = I[g]. \end{aligned} \quad (6.57)$$

where in the last line g is defined by

$$g(t, x, y) = \begin{cases} g_1(t, x, y) & (y \geq 0), \\ g_2(t, x, -y) & (y \leq 0). \end{cases} \quad (6.58)$$

Therefore the last line is the winding number of the configuration g , which is an integer. This was what we wanted to show.

This means that $\Gamma_{\text{WZW}}[g]$ is effectively a two-dimensional integral and can be used in the action of two-dimensional field theories. Indeed, consider the path integral for the partition function of a 2d theory given by

$$Z = \int [Dg] e^{-S[g] + 2\pi i k \Gamma[g]}. \quad (6.59)$$

$\Gamma[g]$ is well-defined modulo 1. Therefore, $e^{2\pi i k \Gamma[g]}$ is well-defined if k is an integer. This is the Wess-Zumino-Witten term in 2d.

Now we can generalize the discussion to 4d. The 5d version of the winding number is known to be given by

$$\mathbb{Z} \ni I[g] = \frac{1}{480\pi^3} \int d^5x \epsilon_{\mu\nu\rho\sigma\tau} \text{tr}(g\partial_\mu g^{-1})(g\partial_\nu g^{-1})(g\partial_\rho g^{-1})(g\partial_\sigma g^{-1})(g\partial_\tau g^{-1}). \quad (6.60)$$

We then consider

$$\Gamma_{\text{wzw}}[g] = \frac{1}{480\pi^3} \int d^5x_{x^5 \geq 0} \epsilon_{\mu\nu\rho\sigma\tau} \text{tr}(g\partial_\mu g^{-1})(g\partial_\nu g^{-1})(g\partial_\rho g^{-1})(g\partial_\sigma g^{-1})(g\partial_\tau g^{-1}). \quad (6.61)$$

Exactly as before, this quantity is well-defined modulo 1, and depends only on the value of g at $x^5 = 0$, i.e. on the value of g on the 4d spacetime. Now we can add to the contribution $2\pi i k \Gamma_{\text{wzw}}[g]$ to the action, which is consistent if k is an integer.

Let us come back to the question of the chiral Lagrangian. We needed a term which involves an epsilon symbol; the Wess-Zumino-Witten term is a good candidate. The coefficient should be of the form $2\pi i k$ for an integer k . There are only two integers in the QCD, N and N_f . N_f already appears in the chiral Lagrangian as the size of the matrix U . Then it sounds reasonable to take $k = N$. This motivates us to upgrade the chiral Lagrangian to the form:

$$S = \int d^4x \left[\frac{F_\pi^2}{4} \text{tr}(\partial_\mu U \partial^\mu U^\dagger) + a \text{tr}(\partial_\mu U \partial^\mu U^\dagger)^2 + \dots \right] + 2\pi N i \Gamma_{\text{wzw}}[U]. \quad (6.62)$$

To leading order, $U \partial_\mu U^{-1} \sim (i/F_\pi) \partial_\mu \pi$ where we parameterize $U = \exp((i/F_\pi)\pi)$. Then,

$$\begin{aligned} 2\pi N i \Gamma_{\text{wzw}}[U] &\sim \frac{N}{240\pi^2 F_\pi^5} \int d^5x_{x^5 \geq 0} \epsilon_{\mu\nu\rho\sigma\tau} \text{tr} \partial_\mu \pi \partial_\nu \pi \partial_\rho \pi \partial_\sigma \pi \partial_\tau \pi \\ &= \frac{N}{240\pi^2 F_\pi^5} \int d^4x \epsilon_{\nu\rho\sigma\tau} \text{tr} \pi \partial_\nu \pi \partial_\rho \pi \partial_\sigma \pi \partial_\tau \pi. \end{aligned} \quad (6.63)$$

This means that the Wess-Zumino-Witten term predicts a five-point coupling among the pions, involving an epsilon tensor, with a striking property that its strength N is quantized to be an integer.

This term also explains why the Skyrmion is a fermion when N is odd. The point is to consider a skyrmion configuration at $t = 0$, and perform an adiabatic rotation by the angle 2π . This gives a specific configuration of U defined on \mathbb{R}^4 , and one simply evaluates $\Gamma_{\text{wzw}}[U]$ on this specific configuration. It turns out that the concrete computation gives $\Gamma_{\text{wzw}}[U] = 1/2$, and therefore

$$e^{2\pi N i \Gamma_{\text{wzw}}[U]} = (-1)^N. \quad (6.64)$$

This means that a Skyrmion produces the sign -1 when rotated by 2π when N is odd, meaning that a Skyrmion is a fermion when N is odd.

Exercise 6.5. Perform this computation, following [Wit83b].

The Wess-Zumino-Witten term is also necessary to reproduce the anomaly of QCD under $SU(N_f) \times SU(\tilde{N}_f)$ in the chiral Lagrangian. This can be done by coupling the gauge fields A and \tilde{A} for the chiral flavor symmetry $SU(N_f) \times SU(\tilde{N}_f)$.¹⁶

7 Renormalizability, effective field theory, and UV completeness

In this final section, let us discuss how to think about the (non)-renormalizability of the Lagrangian of a quantum field theory.

7.1 Assigning dimensions to operators

We use natural unit systems. Every quantity then has an associated dimension. A quantity X is said to have dimension d if X/E^d is dimensionless, where E is some energy scale.¹⁷ This is often denoted by writing $[X] = d$.

The action S of a system appears in the exponent of the integrand of the path integral e^{iS} , and therefore dimensionless. Writing $S = \int d^d x L$, we see $[S] = 0$ and $[d^d x] = -d$, meaning that $[L] = d$. From this we conclude:

- A scalar field with the standard kinetic term $L = (1/2)\partial_\mu\phi\partial^\mu\phi$ then has $[\phi] = d/2 - 1$,
- similarly the gauge field has $[A] = d/2 - 1$,
- while a fermion field with the standard kinetic term $L = \bar{\Psi}\gamma^\mu\partial_\mu\Psi$ has $[\Psi] = (d - 1)/2$.

We can then assign any operator a dimension. For example, a Yukawa interaction $\sim \phi\psi\psi$ has dimension $(3/2)d - 2$. An operator O can appear in the Lagrangian in the form $L \supset \lambda O$, where λ is the coefficient. We see $[O] + [\lambda] = d$.

Suppose we study the system at the energy scale E . The dimensionless quantity characterizing the effect of the term $L \supset \lambda O$ at the scale E is then $\lambda/E^{d-[O]}$. This crude argument already tells us that an operator has a totally different effect in the low-energy physics depending on $[O] < d$, $[O] = d$, and $[O] > d$. We have special terminology for them:

$$\begin{array}{l|l|l} [O] < d & \text{relevant} & \text{(super-)renormalizable} \\ [O] = d & \text{marginal} & \text{renormalizable} \\ [O] > d & \text{irrelevant} & \text{non-renormalizable} \end{array} \quad (7.1)$$

¹⁶The WZW term with A and \tilde{A} given in Witten's original paper [Wit83a] had many typos. This is quite rare for Witten's papers in my experience; usually his papers are trustworthy even in their details. In [Wit83a] he did not provide a way to derive the WZW term with A and \tilde{A} other than saying that it can be found somewhat tediously by trial and error. Systematic methods to obtain it was explained in several papers at the same time, see e.g. [KRS84, CGWS84, KT84, Man85].

¹⁷This definition of dimension is sometimes referred to as the engineering dimension of a quantity in hep-th, when there is a need to distinguish it from the scaling dimension of a quantity, which is about how it behaves under the renormalization group flow. At our rough level of discussion in this section, we do not need to distinguish the two.

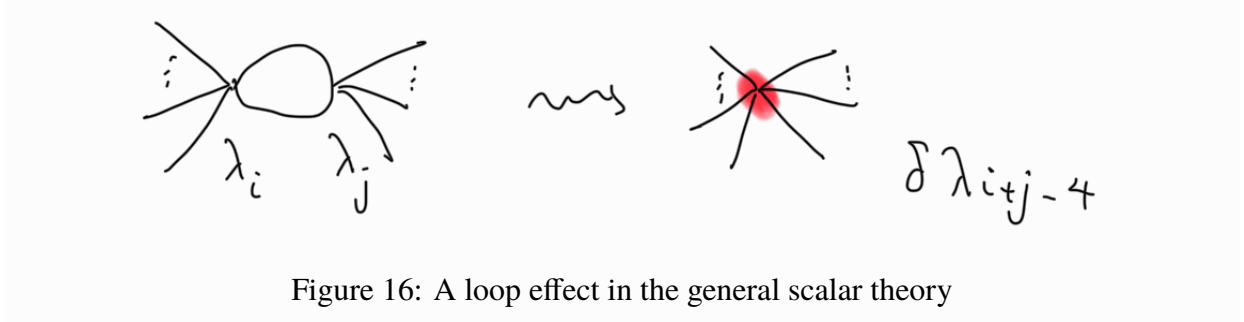


Figure 16: A loop effect in the general scalar theory

where the middle column is used when one emphasizes its effect in the renormalization group flow toward low energy and the right column is used when one emphasizes its effect in the renormalizability in the traditional sense.

Note that there is only a finite number $n(\Delta)$ of operators with dimension below any number Δ ; in particular, there is only a finite number of marginal or relevant operators. But there is an infinite number of irrelevant operators.

The adjectives in the middle column can already be understood by considering how $\lambda/E^{d-[O]}$ behaves when $E \rightarrow 0$. Our aim next is to explain the right column.

7.2 Renormalizability in the traditional sense

To get the idea, let us revisit the computation in Sec. 5.1, in a more general context. Consider the action

$$S = \int d^4x \frac{1}{2} \partial_\mu \phi \partial^\mu \phi + \sum \lambda_i \phi^i. \quad (7.2)$$

As we saw before, $[\phi] = 1$ in $d = 4$, and therefore $[\lambda_i] = 4 - i$.

Let us consider the contribution of the diagram shown in Fig. 16. Note that the combined effect has $i + j - 4$ legs, and therefore gives a quantum correction to λ_{i+j-4} . Very roughly, it is given by

$$\delta\lambda_{i+j-4} \sim \lambda_i \lambda_j \int^\Lambda d^4p \frac{1}{p^2} \frac{1}{p^2} \sim \lambda_i \lambda_j \log \Lambda \quad (7.3)$$

where Λ is the ultraviolet cutoff. Note also that the engineering dimensions of the right hand side and of the left hand side are consistent, since $[\lambda_{i+j-4}] = [\lambda_i] + [\lambda_j] + [\log \Lambda]$. This is infinite as we remove the cutoff $\Lambda \rightarrow \infty$. We need to add $\delta\lambda_{i+j-4}^{\text{by hand}}(\Lambda)$ in the original Lagrangian to cancel this. Such a term is called a counterterm.

Next consider the contribution of the diagram shown in Fig. 17. Similarly we have

$$\delta\lambda_{i-2} \sim \lambda_i \int^\Lambda \frac{d^4p}{p^2} \sim \lambda_i \Lambda^2. \quad (7.4)$$

Again we can check $[\delta\lambda_{i-2}] = [\lambda_i] + [\Lambda^2]$. This correction is infinite as we remove the cutoff, and therefore we need to add the counterterm $\delta\lambda_{i-2}^{\text{by hand}}(\Lambda)$ to cancel it.

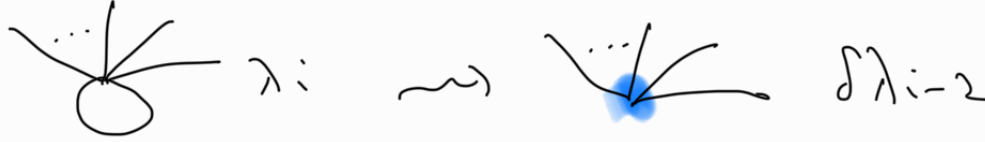


Figure 17: Another loop effect in the general scalar theory

In general, any perturbative computation gives

$$\delta\lambda_{\mathcal{O}} = (\text{a polynomial of } \lambda_{\mathcal{O}'} \text{ and } \Lambda^n \text{ or } \log \Lambda). \quad (7.5)$$

when the cutoff is taken to $\Lambda \rightarrow \infty$. Two obvious remarks are in order:

- Engineering dimension of both sides of the equation should agree.
- The power of Λ is non-negative.

From these observations we find the following immediate conclusion: If the original Lagrangian contains only operators \mathcal{O} with $[\lambda_{\mathcal{O}}] \geq 0$, the perturbative divergences only appear in $\delta\lambda_{\mathcal{O}'}$ with $[\lambda_{\mathcal{O}'}] \geq 0$. In other words, we have

If the original Lagrangian contains (super-)renormalizable operators, one needs to add counterterms to (super-)renormalizable operators. In particular, one only needs a finite number of counterterms in this case. Such a theory is called *renormalizable* in the traditional sense.

In contrast, if one has an operator with $[\lambda_{\mathcal{O}}] < 0$, simply by using this operator many times in a diagram, we can generate divergences in $\delta\lambda_{\mathcal{O}'}$ with arbitrarily high dimensions. Paraphrasing, we have

If the original Lagrangian contains even a single non-renormalizable operators, one needs to specify infinitely many counterterms. Such a theory is called *non-renormalizable* in the traditional sense.

7.3 Non-renormalizable theories as effective theories

Theories which are non-renormalizable in the traditional sense was not liked in the olden days, because one needs to specify infinite number of counterterms by hand in the process of the computations. This naively seemed to remove any predictability from the theory: if we have an infinite number of knobs to tune, we would be able to tune anything. But this is not the case as the following example shows.

Consider a theory of fermions with the following action:

$$S = \int d^4x [\bar{\Psi}\gamma^\mu\partial_\mu\Psi + G(\bar{\Psi}\Psi)^2 + \dots], \quad (7.6)$$

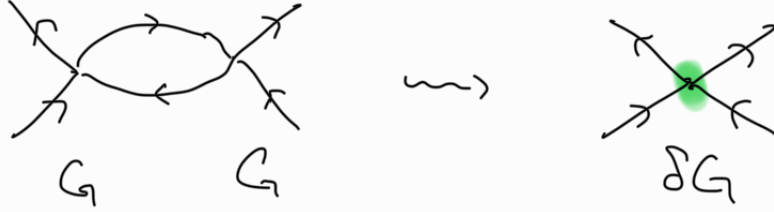


Figure 18: A loop effect of the four-Fermi interactions

where $[G] = -2$ and therefore this interaction is irrelevant=non-renormalizable. This form of the interaction is known as the 4-Fermi interaction.

Now consider the effect of the diagram shown in Fig. 18. This is quadratically divergent, since it is of order $G^2 \int d^4k (1/k)^2$. Then we have, by dimensional analysis, that

$$G + \delta G(P) \sim G + G^2(c\Lambda^2 + c'P^2 + c''P^2 \log(\Lambda/P^2)) + \dots \quad (7.7)$$

where c etc. are dimensionless constants which can be determined by explicit computations, and P is the typical momentum scale of the diagram. Adding counterterms remove the dependence on Λ , but we still have a logarithmic dependence on the momentum scale P

$$G(P) \sim G(P_0) + c''G(P_0)^2 P^2 \log(P^2/P_0^2) + \dots \quad (7.8)$$

with a computable coefficient c'' .

Note that after the Fourier transformation to the position space, correction terms polynomial in P^2 lead to delta-function potentials or its derivatives, whereas logarithmic terms can lead to potentials of the form $1/r^3$. To see this, recall that the massless exchange $1/p^2$ leads to a Coulombic force since

$$\int d^3p e^{ixp} \frac{1}{p^2} \sim \frac{1}{r}, \quad r = |x|. \quad (7.9)$$

Then the correction of the order $(p^2)^n$ leads to

$$\int d^3p e^{ixp} (p^2)^n \sim (\partial^2)^n \delta(x). \quad (7.10)$$

In contrast, the correction of the order $\log p^2$ corresponds to

$$\int d^3p e^{ixp} \log(p^2) \sim \frac{1}{r^3}. \quad (7.11)$$

In this sense the terms logarithmic in the momentum p can directly affect the long-range potential.

Coming back to the general discussion, the point is that any local term in the Lagrangian is built from fields and derivatives. Therefore all that a local counterterm can directly give is a polynomial function in P , and a local term cannot produce logarithmic terms unless via loops. Therefore, the

coefficient of the logarithmic term is a prediction which cannot be tuned by the infinite number of knobs which are the coefficient of the non-renormalizable terms.

That said, the presence of the non-zero coefficient G for the non-renormalizable interaction such as $(\bar{\Psi}\Psi)^2$ means that the perturbative computation in that theory stops being meaningful when the energy scale E is so high that the dimensionless coefficient GE^2 is of order 1. This theory stops being meaningful (or effective) at this energy scale.

A quantum field theory which makes sense only below a certain energy scale is called an *effective field theory*.¹⁸ whereas a quantum field theory which makes sense at an arbitrarily small scale is called a UV complete theory.¹⁹

From our discussions above we see

$$\{\text{non-renormalizable theories}\} \subset \{\text{effective field theories}\} \quad (7.12)$$

but not all renormalizable theories are UV complete. For example, the ϕ^4 theory we discussed in Sec. 5.1 and the QED in Sec. 5.2.2 are both renormalizable in the sense that the Lagrangian only contains renormalizable terms. But we already saw there that there is a Landau pole in both cases, and the theories make sense only up to some energy scale. Therefore they are not UV complete.

The QCD we discussed in Sec. 5.2.3 is renormalizable for any N and N_F . It has a Landau pole when $N_F/N > 11/2$, which means that it is an effective field theory. When $N_F/N < 11/2$, the system is asymptotically free, and therefore it is a UV complete theory.

7.4 ‘Completion’ of an effective theory

What happens at this scale when the effective field theory breaks down? One scenario is that there is an interesting strongly-coupled dynamics which cannot be analyzed by perturbation theory. This scenario is known to be realized in examples.

Another less drastic scenario is that there is a heavy particle whose mass M is around this scale, i.e. so that $GM^2 \sim 1$, which needs to be taken into account. The 4-Fermi model was first introduced to describe the decay of a neutron to a proton, an electron, and an anti-neutrino; this is now known to be mediated by a W-boson. A simpler toy example is the following.

7.4.1 4-Fermi theory

Consider the theory of a scalar and a fermion with the Lagrangian

$$S = \int d^4x \left[\frac{1}{2} \partial_\mu \phi \partial^\mu \phi + \frac{1}{2} M^2 \phi^2 + \bar{\Psi} \not{\partial} \Psi + y \phi \bar{\Psi} \Psi \right] \quad (7.13)$$

Note that the Yukawa interaction has $[y] = 0$ and therefore is marginal; this model is renormalizable. We now consider the tree level diagram shown in Fig. 19. We see that it generates the

¹⁸I do not think that the adjective ‘effective’ is not used here in the sense that “wow, this theory is very effective!”. It is used rather in the sense that it works although not formally completely legit.

¹⁹Whether people mean a UV complete QFT or an effective QFT by the terminology ‘QFT’ without the adjective depends on the subfield of physics and/or mathematical physics. One needs some caution, therefore, when one reads textbooks and/or review articles on the general properties of QFTs from different subfields one is in.

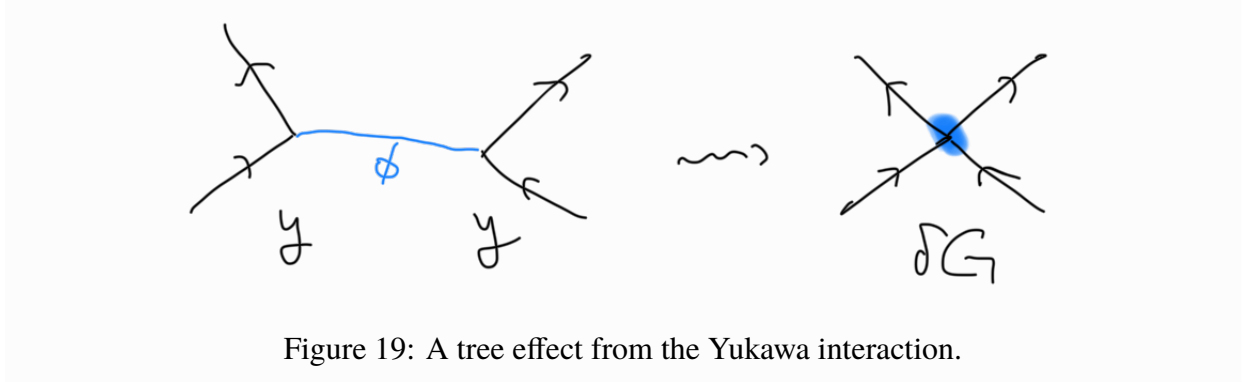


Figure 19: A tree effect from the Yukawa interaction.

4-Fermi interaction of the strength

$$\delta G \sim \frac{y^2}{p^2 + M^2} \quad (7.14)$$

$$\rightarrow \frac{y^2}{M^2} \quad (7.15)$$

where we took the low energy limit $|p| \lll M$. This operation of replacing the effect of a massive particle by newly generated interaction vertices in a Lagrangian is called to *integrate out the field* ϕ . This terminology is based on the path integral formalism, where we consider performing the path integral of the field ϕ first.

This means that the 4-fermi theory (7.6) with $G \sim y^2/M^2$ can be thought of as the low-energy effective description of the system (7.13) of a massive scalar and a fermion. If one does not have enough energy of order M to create the ϕ particle, one cannot see the difference between the two systems. But the fact that the 4-fermi theory (7.6) contains a non-renormalizable interaction means that something *must* happen at a certain high energy.

7.4.2 Chiral Lagrangian

The chiral Lagrangian of pions we discussed to some extent in Sec. 6.4 is also an example of an effective field theory. In the form

$$S = \int d^4x \frac{F_\pi^2}{4} \text{tr} \partial_\mu U \partial^\mu U^\dagger \quad (7.16)$$

it might look like it contains only the kinetic term, but in terms of the pion fields as in (6.43) and (6.44), the Lagrangian contains infinitely many terms with arbitrarily high dimensions. This means that this Lagrangian stops being meaningful or effective when the energy is too high. Indeed, we now know that the system in the short distance limit is described by QCD, whose Lagrangian was given in (6.2), which can be checked to be renormalizable.

7.4.3 Perturbative unitarization

Here we note that the Yukawa interaction and the chiral Lagrangian were completed in a rather different manner: The Yukawa interaction was saved by simply introducing another massive field, which when integrated out would generate the non-renormalizable interaction. The chiral Lagrangian was saved rather by using a totally different Lagrangian of QCD in the ultraviolet.

The former method of UV completion is called the perturbative unitarization. The terminology comes from the following. One effect of the non-renormalizable interaction is a too-fast growth of the scattering amplitude in the high energy region, which at face value would correspond to a scattering probability exceeding one, thus destroying the unitarity. Adding a massive particle removes this problem, thus making the system unitary perturbatively.

Another famous example of perturbative unitarization is the following. Consider the 4-fermi interaction

$$\sim G_F (\bar{\Psi}_p \gamma^\mu \Psi_n) (\bar{\Psi}_e \gamma_\mu \Psi_\nu) \quad (7.17)$$

describing the weak-decay process. (More precisely we need to include the chirality projector.) This is non-renormalizable, and requires a UV completion. This can be achieved by introducing a massive charged vector boson W_μ^\pm , just as the interaction $G(\bar{\Psi}\Psi)^2$ can be generated by the exchange of a massive scalar.

Now, the interaction of the massive vector bosons themselves is known to lead to the violation of perturbative unitarity in the high energy region. This is can be saved by introducing the Higgs field, which generates the mass of the vector boson by the Higgs mechanism. If the Higgs boson is too massive, the perturbative unitarity is violated before the Higgs boson comes to save the day. In this manner one can predict the upper bound of the Higgs boson mass to be around 1 TeV [LQT77a, LQT77b].

7.5 Standard Model as an effective field theory

Let us next consider the Standard Model we discussed in Sec. 2. The terms in the Lagrangian we introduced there were all (super-)renormalizable. Therefore the Standard Model is renormalizable. However, we know that there are more to our world than is contained in this Standard Model, and therefore the Standard Model is an effective field theory. Even purely theoretically, the $U(1)$ coupling has a Landau pole, so the Standard Model is not UV complete, and can only be an effective theory.

This means that one needs to study the Standard Model as an effective field theory. The first step should then be to enumerate all possible operators one can potentially add to the Lagrangian, from low dimensions to high dimensions. The importance was of course noticed from the early days, and the task was carried out at dimension 5 by [Wei80], dimension 6 by [BW86, GIMR10, AJMT13], dimension 7 and 8 by [Leh14, LM15] in an ad hoc manner. Somewhat surprisingly, a systematic method to perform this computation was only introduced in [HLMM15, HLMM17]. One complication is that as a term in the Lagrangian, one needs to remove a total derivative, since

$$L = \int d^4x \partial_\mu O^\mu \quad (7.18)$$

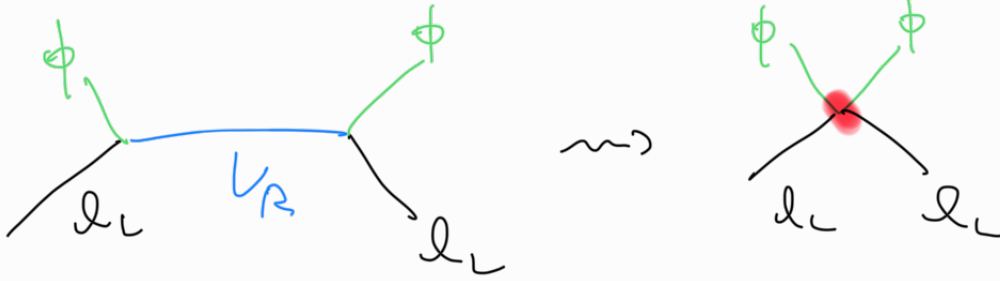


Figure 20: Mass terms for $\nu_L = (\phi\ell_L)$ generated from the exchange of ν_R .

can be integrated by parts and vanishes under a usual choice of boundary conditions at infinity.

Exercise 7.1. Count how many non-renormalizable operators there are in the Standard Model, by going through these papers.

At this point it is important to mention the following. In the Lagrangian of the standard model we wrote in Sec. 2, we included in it the right-handed neutrino fields $\bar{\nu}_R^i$ and their Majorana mass term

$$M_{ij}\bar{\nu}_R^i\nu_R^j + c.c. \quad (7.19)$$

This is definitely not the conventional way of saying things. Usually the standard model Lagrangian refers to the state of affairs before the nonzero neutrino mass for ν_L was experimentally discovered 20 years ago. Without $\bar{\nu}_R$, one cannot write mass terms for the neutrino if we only allow (super)-renormalizable terms in the Lagrangian. The neutrino is then automatically massless.

Exercise 7.2. Check this.

One way to have nonzero neutrino mass for ν_L without having $\bar{\nu}_R$ is to allow for non-renormalizable terms. Indeed, one can integrate out $\bar{\nu}_R$ as in Fig. 20, leading to the operator of the form

$$(\ell_L)_\alpha^u\phi_u(\ell_L)_\beta^v\phi_v\epsilon^{\alpha\beta} + c.c. \quad (7.20)$$

whose coefficient is of order $(Y^{\text{neutrino}})^2/M_{\text{Majorana}}$. This operator, after ϕ is replaced by the vev, generates the neutrino mass terms. Note that this operator is of dimension 5; in fact this operator (together with the complex conjugate) is the only dimension 5 operator in the standard model. The whole mechanism is known as the see-saw mechanism, since M_{Majorana} of $\bar{\nu}_R$ appears in the denominator of the operator for the mass of ν_L .

We also note that the renormalizable standard-model Lagrangian without $\bar{\nu}_R$ automatically has the baryon number symmetry and the lepton number symmetry classically, where we assign baryon charge $B = +1/3$ to Q , \bar{u}_R and \bar{d}_R , and lepton charge $L = +1$ to ℓ_L and \bar{e}_R . Because of the anomaly, only the combination $B - L$ is a quantum mechanical symmetry.

Exercise 7.3. Check this.

The Majorana mass term $M\nu_R\nu_R$ breaks this $B-L$ symmetry, and therefore the renormalizable Lagrangian with $\overline{\nu_R}$ does not have this $B-L$ symmetry unless the Majorana mass term is zero. The operator (7.20) also breaks this $B-L$ symmetry.

This means that all of the baryon-number or lepton-number violating terms in the standard model Lagrangian is non-renormalizable or irrelevant, which therefore carry naturally small coefficients. This explains the experimental fact that the baryon and the lepton number are very well conserved.

Exercise 7.4. The process which violates individual baryon and lepton number symmetry but keeps the $B-L$ symmetry exists in nature but is known to be a very small effect in practice. Study about it.

7.6 Renormalizability of Yang-Mills and gravity in various dimensions

Let us next consider the renormalizability of Yang-Mills and gravity in various dimensions. The Yang-Mills theory has the Lagrangian

$$S_{\text{YM}} = \int d^d x \frac{1}{2g^2} \text{tr} F_{\mu\nu} F^{\mu\nu} \quad (7.21)$$

where we used the geometric convention

$$F_{\mu\nu} = i[D_\mu, D_\nu], \quad D_\mu = \partial_\mu + A_\mu. \quad (7.22)$$

This forces $[A] = 1$ independent of the spacetime dimension d and put the dimensionality to the coupling constant, $[g] = (4-d)/2$.

If we expand the Lagrangian and rescale A_μ to have the canonical dimension, one finds that

$$S_{\text{YM}} \sim \int d^d x [\partial A \partial A + g \partial A A A + g^2 A A A A + \dots] \quad (7.23)$$

which means that the dimension of g is what controls the renormalizability.²⁰

This means that the Yang-Mills theory is super-renormalizable when $d < 4$, renormalizable when $d = 4$, and non-renormalizable when $d > 4$. The dimension 4 is sometimes referred to as the critical dimension of the gauge theory.

As for the gravity, the Einstein-Hilbert action is

$$S_{\text{gravity}} = \int d^d x \frac{1}{16\pi G_N} \sqrt{g} R. \quad (7.24)$$

²⁰In the renormalizability of a gauge theory, one needs to make sure that the possible divergences can be absorbed by gauge invariant counterterms, in order not to spoil the gauge invariance. This point was beautifully treated in Hori-san's installment of the QFT class [Hor18], which is highly recommended.

d	1	2	3	4	5	6	...
Yang-Mills	✓✓	✓✓	✓✓	✓	△	△	
gravity	✓✓	✓	△	△	△	△	

Table 1: Renormalizability of Yang-Mills and gravity. ✓✓: super-renormalizable, ✓: renormalizable, △: non-renormalizable.

Our convention is again geometric; the metric determines the proper length via

$$ds^2 = g_{\mu\nu} dx^\mu dx^\nu, \quad (7.25)$$

and therefore $[g_{\mu\nu}] = 0$. Since $R \sim \partial\partial g$, we see that $[G_N] = 2 - d$. To perform perturbation theory, we assume $g_{\mu\nu}$ is close to a flat metric

$$g_{\mu\nu} = \eta_{\mu\nu} + \sqrt{G_N} h_{\mu\nu}. \quad (7.26)$$

Then the action becomes

$$S_{\text{gravity}} \sim \int d^d x [\partial h \partial h + \sqrt{G_N} \partial^2 h^3 + \dots]. \quad (7.27)$$

Therefore we see again that the renormalizability is controlled by the dimension of G_N . We find that the gravity is super-renormalizable when $d < 2$, renormalizable when $d = 2$, and non-renormalizable when $d > 2$. The critical dimension of gravity is $d = 2$. The situation concerning renormalizability of Yang-Mills and gravity is summarized in Table 1.

It is worth emphasizing that quantum gravity in $d = 1$ and in $d = 2$ is rather well-understood using the standard methods of quantum field theory, to the same degree that we understand quantum Yang-Mills gauge theory in $d = 3$ and $d = 4$. It is true that we do not yet know the UV completion of the quantum gravity in $d = 4$, but the situation does not a priori require a conceptual revolution, in contrast to what one would often read and hear in expository articles to general public, by string theories or loop quantum gravity theorists.

One example illustrating this fact is that the logarithmic dependence on the momentum of the gravitational coupling by the one-loop effects of graviton can be reliably computed, just as in the case of the 4-Fermi theory we discussed in Sec. 7.3. By performing a Fourier transformation to translate it into the gravitational potential, one finds [BBDH02]

$$V(r) = -\frac{Gm_1 m_2}{r} \left[1 + 3 \frac{G(m_1 + m_2)}{r} + \frac{41}{10\pi} \frac{G}{r^2} \right] + \dots \quad (7.28)$$

where the first term is the Newton potential, the second term is the classical general relativistic correction, and the third term is the one-loop effect. For more on the quantum gravity as an effective field theory, see the excellent review article [Bur03].

As a final comment, I would like to mention the study of the perturbative unitarization of 4d gravity by T.-C. Huang, Y.-T. Huang, and N. Arkani-Hamed [Hua16, AH16]. As 4d gravity is

non-renormalizable, one might want to find a UV completion. One way would be a perturbative unitarization, as in the case of the 4-Fermi interaction.

Recall that in the case of the 4-fermi interaction $\bar{\Psi}\gamma^\mu\Psi\bar{\Psi}\gamma_\mu\Psi$, it was a two-step process: one first introduces a massive charged vector boson, and then one introduces a Higgs boson. In the case of 4d gravity, they found that this process continues indefinitely: denoting the mass squared of the first massive particle we add by M^2 , they found that perturbative unitarization requires an addition of an infinite series of massive particles of mass squared, given by

$$m_n^2 = Nm^2. \tag{7.29}$$

This spectrum is a typical one for a perturbative string theory.

String theory is usually introduced by fiat by the assumption that things were made of strings. The argument of [Hua16, AH16] shows that it arises naturally by trying to unitarize 4d gravity perturbatively.

References

- [AC79] Y. Aharonov and A. Casher, *Ground state of a spin- $\frac{1}{2}$ charged particle in a two-dimensional magnetic field*, *Phys. Rev. A* **19** (1979) 2461–2462.
- [AH16] N. Arkani-Hamed, *Towards Deriving String Theory as the Weakly Coupled UV Completion of Gravity*. <https://member.ipmu.jp/yuji.tachikawa/stringsmirrors/2016/main/Nima%208.3.pdf>. Talk at Strings 2016, Beijing.
- [AJMT13] R. Alonso, E. E. Jenkins, A. V. Manohar, and M. Trott, *Renormalization Group Evolution of the Standard Model Dimension Six Operators Iii: Gauge Coupling Dependence and Phenomenology*, *JHEP* **04** (2014) 159, [arXiv:1312.2014](https://arxiv.org/abs/1312.2014) [hep-ph].
- [BBDH02] N. E. J. Bjerrum-Bohr, J. F. Donoghue, and B. R. Holstein, *Quantum Gravitational Corrections to the Nonrelativistic Scattering Potential of Two Masses*, *Phys. Rev. D* **67** (2003) 084033, [arXiv:hep-th/0211072](https://arxiv.org/abs/hep-th/0211072). [Erratum: *Phys. Rev. D* **71**, 069903(2005)].
- [BM74] A. A. Belavin and A. A. Migdal, *Calculation of Anomalous Dimensions in Non-Abelian Gauge Field Theories*, *Pisma Zh. Eksp. Teor. Fiz.* **19** (1974) 317–320.
- [Bur03] C. P. Burgess, *Quantum Gravity in Everyday Life: General Relativity as an Effective Field Theory*, *Living Rev. Rel.* **7** (2004) 5–56, [arXiv:gr-qc/0311082](https://arxiv.org/abs/gr-qc/0311082).
- [BW86] W. Buchmuller and D. Wyler, *Effective Lagrangian Analysis of New Interactions and Flavor Conservation*, *Nucl. Phys.* **B268** (1986) 621–653.
- [BZ82] T. Banks and A. Zaks, *On the Phase Structure of Vector-Like Gauge Theories with Massless Fermions*, *Nucl. Phys.* **B196** (1982) 189–204.

- [Cas74] W. E. Caswell, *Asymptotic Behavior of Nonabelian Gauge Theories to Two Loop Order*, *Phys. Rev. Lett.* **33** (1974) 244.
- [CGWS84] K.-c. Chou, H.-y. Guo, K. Wu, and X.-c. Song, *On the Gauge Invariance and Anomaly Free Condition of Wess-Zumino-Witten Effective Action*, *Phys. Lett.* **134B** (1984) 67–69.
- [DeG15] T. DeGrand, *Lattice Tests of Beyond Standard Model Dynamics*, *Rev. Mod. Phys.* **88** (2016) 015001, [arXiv:1510.05018](https://arxiv.org/abs/1510.05018) [hep-ph].
- [GIMR10] B. Grzadkowski, M. Iskrzynski, M. Misiak, and J. Rosiek, *Dimension-Six Terms in the Standard Model Lagrangian*, *JHEP* **10** (2010) 085, [arXiv:1008.4884](https://arxiv.org/abs/1008.4884) [hep-ph].
- [HLMM15] B. Henning, X. Lu, T. Melia, and H. Murayama, 2, 84, 30, 993, 560, 15456, 11962, 261485, ...: *Higher Dimension Operators in the SM EFT*, *JHEP* **08** (2017) 016, [arXiv:1512.03433](https://arxiv.org/abs/1512.03433) [hep-ph].
- [HLMM17] ———, *Operator bases, S-matrices, and their partition functions*, *JHEP* **10** (2017) 199, [arXiv:1706.08520](https://arxiv.org/abs/1706.08520) [hep-th].
- [Hor18] K. Hori, *Elementary Particle Theory*. <https://member.ipmu.jp/kentaro.hori/Courses/EPP/>.
- [Hua16] Y. T. Huang, *Lessons from perturbative unitarity in graviton scattering amplitudes*. <http://www2.yukawa.kyoto-u.ac.jp/~qft.web/2016/slides/huang.pdf>. Talk at Strings and Fields 2016, YITP, Kyoto U.
- [IS15] K. Intriligator and F. Sannino, *Supersymmetric Asymptotic Safety is Not Guaranteed*, *JHEP* **11** (2015) 023, [arXiv:1508.07411](https://arxiv.org/abs/1508.07411) [hep-th].
- [Jon74] D. R. T. Jones, *Two Loop Diagrams in Yang-Mills Theory*, *Nucl. Phys.* **B75** (1974) 531.
- [JW06] A. Jaffe and E. Witten, *Quantum Yang-Mills theory*, The millennium prize problems, Clay Math. Inst., Cambridge, MA, 2006, pp. 129–152. <https://www.claymath.org/millennium-problems/yang%E2%80%93mills-and-mass-gap>.
- [Kra03] H. Kragh, *Magic number: A partial history of the fine-structure constant*, *Archive for History of Exact Sciences* **57** (2003) 395–431.
- [KRS84] O. Kaymakcalan, S. Rajeev, and J. Schechter, *Nonabelian Anomaly and Vector Meson Decays*, *Phys. Rev.* **D30** (1984) 594.
- [KT84] H. Kawai and S. H. H. Tye, *Chiral Anomalies, Effective Lagrangian and Differential Geometry*, *Phys. Lett.* **140B** (1984) 403–407.

- [Lam97] S. K. Lamoreaux, *Demonstration of the casimir force in the 0.6 to $6\mu\text{m}$ range*, *Phys. Rev. Lett.* **78** (1997) 5–8.
- [Leh14] L. Lehman, *Extending the Standard Model Effective Field Theory with the Complete Set of Dimension-7 Operators*, *Phys. Rev.* **D90** (2014) 125023, arXiv:1410.4193 [hep-ph].
- [LM15] L. Lehman and A. Martin, *Low-Derivative Operators of the Standard Model Effective Field Theory via Hilbert Series Methods*, *JHEP* **02** (2016) 081, arXiv:1510.00372 [hep-ph].
- [LMS15] D. F. Litim, M. Mojaza, and F. Sannino, *Vacuum Stability of Asymptotically Safe Gauge-Yukawa Theories*, *JHEP* **01** (2016) 081, arXiv:1501.03061 [hep-th].
- [LQT77a] B. W. Lee, C. Quigg, and H. B. Thacker, *The Strength of Weak Interactions at Very High-Energies and the Higgs Boson Mass*, *Phys. Rev. Lett.* **38** (1977) 883–885.
- [LQT77b] ———, *Weak Interactions at Very High-Energies: the Role of the Higgs Boson Mass*, *Phys. Rev.* **D16** (1977) 1519.
- [LS14] D. F. Litim and F. Sannino, *Asymptotic Safety Guaranteed*, *JHEP* **12** (2014) 178, arXiv:1406.2337 [hep-th].
- [Man85] J. L. Manes, *Differential Geometric Construction of the Gauged Wess-Zumino Action*, *Nucl. Phys.* **B250** (1985) 369–384.
- [tH76] G. 't Hooft, *Symmetry Breaking Through Bell-Jackiw Anomalies*, *Phys. Rev. Lett.* **37** (1976) 8–11. [,226(1976)].
- [tH86] ———, *How Instantons Solve the U(1) Problem*, *Phys. Rept.* **142** (1986) 357–387.
- [Wei80] S. Weinberg, *Varieties of Baryon and Lepton Nonconservation*, *Phys. Rev.* **D22** (1980) 1694.
- [Wit83a] E. Witten, *Current Algebra, Baryons, and Quark Confinement*, *Nucl. Phys.* **B223** (1983) 433–444.
- [Wit83b] ———, *Global Aspects of Current Algebra*, *Nucl. Phys.* **B223** (1983) 422–432.
- [Yan05] C.-N. Yang, *Selected papers (1945–1980)*, World Scientific, 2005.