

Lattice quantum field theory

Lecture in SS 2020 at the KFU Graz

Axel Maas

Contents

1	Introduction	1
2	Scalar fields on the lattice	4
2.1	The path integral	4
2.2	Euclidean space-time lattices	7
2.3	Discrete analysis	10
2.4	The free scalar field	11
2.5	ϕ^4 theory on the lattice	14
2.6	Lattice perturbation theory	15
2.6.1	Particle properties	15
2.6.2	Interactions	17
2.6.3	Hopping expansion	20
2.7	The phase diagram and the continuum limit	22
2.8	Detecting phase transitions	26
2.9	Internal and external symmetries on the lattice	28
2.9.1	External symmetries	28
2.9.1.1	Translation symmetry	29
2.9.1.2	Reflection positivity	29
2.9.1.3	Rotation symmetry and angular momentum	31
2.9.2	Internal symmetries	33
3	Measurements	35
3.1	Expectation values	35
3.2	Spectroscopy and energy levels	36
3.3	Bound states	38
3.4	Resonances, scattering states, and the Lüscher formalism	39
3.5	General correlation functions	44

4	Monte-Carlo simulations	48
4.1	Importance sampling and the Markov chain	48
4.2	Metropolis algorithm	52
4.3	Improving algorithms	54
4.3.1	Acceleration	54
4.3.2	Overrelaxation	54
4.4	Statistical errors	55
4.4.1	Determination and signal-to-noise ratio	55
4.4.2	Autocorrelations, thermalization, and critical slowing down	57
4.5	Systematic errors	60
4.5.1	Lattice artifacts	60
4.5.2	Overlap and reweighting	61
4.5.3	Ergodicity	62
4.6	Spectroscopy	63
4.6.1	Time-slice averaging	63
4.6.2	Smearing	64
4.6.3	Variational analysis	65
4.6.4	Resonances	67
4.7	Phase transitions revisited	68
4.8	Spontaneous symmetry breaking in numerical simulations	69
5	Gauge fields on the lattice	71
5.1	Abelian gauge theories	72
5.2	Non-Abelian gauge theories	76
5.3	Gauge-invariant observables and glueballs	78
5.4	Numerical algorithms for gauge fields	82
5.5	Gauge-fixing	85
5.6	Perturbation theory revisited	90
5.7	Scaling and the continuum limit	92
5.8	The strong-coupling expansion	94
5.8.1	Construction	94
5.8.2	Wilson criterion	96
5.9	Improved actions	98
5.10	Topology	100
5.11	Weak interactions and the Higgs	101

6	Fermions	104
6.1	Naive fermions and the doubling problem	104
6.2	Wilson fermions	106
6.3	Staggered fermions	108
6.4	Ginsparg-Wilson fermions	109
6.5	QCD	111
6.6	Perturbation theory for fermions	113
6.7	Fermionic expectation values	114
6.8	Hadron spectroscopy	117
6.9	Algorithms for fermions	119
6.9.1	Molecular dynamic algorithms	120
6.9.2	Molecular dynamics for gauge fields	122
6.9.3	Updating QCD	123
7	Finite temperature and density	126
7.1	Finite temperature	126
7.2	Finite density and the sign problem	128
	Index	130

Chapter 1

Introduction

Lattice quantum field theories are essentially quantum field theories which are not defined in continuum Minkowski time, but rather in a finite Euclidean volume on a discrete set of points, the lattice¹. It is at first sight very dubious that both formulations should have anything to do with each other, or that this may be useful. But this is not true.

As will be discussed at length, a lattice quantum field theory can, in a very precise sense, be taken as an approximation of an ordinary continuum² field theory. In fact, the lattice can be taken as a regulator of the theory, both in the ultraviolet, due to the discrete nature of the set of points there is a maximum energy given the smallest distance between points, as well as in the infrared, due to the finite volume. Thus, a lattice theory can be taken as a regularized version of the corresponding continuum field theory. Because the number of points is finite, there is only a finite number of degrees of freedom. Therefore, it is an approximation of a quantum field theory by a quantum mechanical system.

Because of this, many statements in lattice field theory can be made on a quite rigorous level, because they have only to cover a finite number of degrees of freedom. There are therefore many powerful and rigorous statements on lattice theories. Whether they also hold for the continuum theory is actually rarely known. Taking the limit of an infinite number of degrees of freedom usually plays havoc with various assumptions. It is therefore usually unknown whether the approximation by a lattice version of the theory is good. Fortunately, for many theories this seems to be the case, though rather than proofs there is usually only circumstantial evidence.

¹Note that a lattice field theory can also be defined in Minkowski space-time, and with various things being infinite rather than finite. However, this is beyond the mainstream lattice to be presented in this lecture.

²Though precisely continuum does say nothing about the different metric nor the finite volume, it has become customary to just use the statement continuum to distinguish between both versions of the theory. This is somewhat an abuse of language.

This problem is somewhat alleviated by the fact that ultimately this question is somewhat academic. Since the universe is (likely) finite and quantum gravity appears at some ultraviolet scale, any theory on flat Minkowski space-time is anyhow only expected to be relevant over a finite distance and up to a finite energy in the sense of a low-energy effective theory. Therefore, the lattice version of a quantum field theory may in the end be actually even be a better approximation of nature than a continuum theory, if both should not coincide.

These remarks show the conceptual importance of lattice quantum field theories. There is also a technical important one. As will be seen, the lattice version of a theory is accessible to numerical evaluations of a very particular kind: It is possible to approximate the path integral for any observable with, in principle, arbitrary precision by a numerical calculation. Especially, in many cases the required computational time grows only like a (small) power of the number of points of the lattice. Therefore, such calculations are actually feasible on nowadays computers. Since these numerical evaluations are, up to certain error sources which can be improved, exact, this implies that all information is sampled, including non-perturbative effects. This possibility makes lattice quantum field theories nowadays to a mainstay tool for non-perturbative calculations in lattice quantum field theory, though there are several areas where the numerical cost is actually still too large. This is probably an even more important reason to have a look a lattice quantum field theories.

All of this will be discussed during this lecture. The particular emphasis is on the techniques and concepts of lattice quantum field theory. The phenomenology of the theories investigated using these methods is not the subject of this lecture, but rather of various other ones. It is thus a lecture about techniques. The advantage is that these techniques can be applied to essentially any quantum field theory as a powerful tool.

The number of books on this topic is somewhat limited. The ones which have been used to prepare this lecture areas

- Böhm et al. “Gauge theories” (Teubner)
- DeGrand et al. “Lattice methods for quantum chromodynamics (World Scientific)
- Lang et al. “Quantum chromodynamics on the lattice” (Springer)
- Montvay et al. “Quantum fields on a lattice” (Cambridge)
- Rothe “Lattice gauge theories” (World Scientific)
- Seiler “Gauge theories as a problem of constructive quantum field theory and statistical mathematics” (Springer)

As the titles in this list shows, one of the central subjects where lattice quantum field theory has shown its most spectacular successes in the past has been quantum chromodynamics. Due to the strength of its interactions, QCD exhibits a large number of hard to control nonperturbative features³. It is here where the ability to deal numerically with the nonperturbative features excelled.

³Note that strong interactions and nonperturbative features are not equivalent. The best example is the rich solid state physics, which is entirely due to non-perturbative interactions of the weakly interacting QED.

Chapter 2

Scalar fields on the lattice

2.1 The path integral

To define a lattice theory the path-integral formulation is the method of choice. Since defining the path integral itself is usually done using a lattice approximation, it is useful to consider this in more detail¹.

Since the path integral formulation is as axiomatic as is canonical quantization, it cannot be deduced. However, it is possible to motivate it.

A heuristic reasoning is the following. Take a quantum mechanical particle which moves in time T from a point a of origin to a point b of measurement. This is not yet making any statement about the path the particle followed. In fact, in quantum mechanics, due to the superposition principle, a-priori no path is preferred. Therefore, the transition amplitude U for this process must be expressible as

$$U(a, b, T) = \sum_{\text{All paths}} e^{i \cdot \text{Phase}} \quad (2.1)$$

which are weighted by a generic phase associated with the path. Since all paths are equal from the quantum mechanical point of view, this phase must be real. Thus it remains only to determine this phase. Based on the correspondence principle, in the classical limit the classical path must be most important. Thus, to reduce interference effects, the phase should be minimal for the classical path. A function which implements this is the classical action S , determined as

$$S[C] = \int_C dt L,$$

¹It is actually possible to define, even mathematically rigorous in the non-interacting case, the path integral directly in the continuum. However, this requires more general ways of integrating, as quantum fields are usually non-continuous functions, which cannot be treated by Riemann integration.

where the integral is over the given path C from a to b , and the action is therefore a functional of the path S and the classical Lagrange function L . Of course, it is always possible to add a constant to the action without altering the result. Rewriting the sum as a functional integral over all paths, this yields already the definition of the functional integral

$$U(a, b, T) = \sum_C e^{iS[C]} \equiv \int \mathcal{D}C e^{iS[C]}.$$

This defines the quantum mechanical path integral.

It then remains to give this functional integral a more constructive meaning, such that it becomes a mathematical description of how to determine this transition amplitude. The most useful approach so far for non-trivial interacting theories is the intermediate use of a lattice, i. e., in fact lattice field theory as described in this lecture. However, even in this case there are still conceptual and practical problems, so that the following remains often an unproven procedure.

Thus, it is useful to first further consider the quantum-mechanical case, to make all of the above better defined.

The starting point was the transition amplitude. In quantum mechanics, this amplitude is given by

$$U(a, b, T = t_N - t_0) = \langle a, t_N | e^{-iHT} | b, t_0 \rangle.$$

In the next step, insert at intermediate times a sum, or integral in cases of a continuous spectrum, over all states

$$U(a, b, T) = \sum_i \langle a, t_N | e^{-iH(t_N-t_1)} | i, t_1 \rangle \langle i, t_1 | e^{-iH(t_1-t_0)} | b, t_0 \rangle.$$

By this, the transition amplitude is expressed by a sum over all possible intermediate states, already quite in the spirit of (2.1). To fully embrace the idea, divide the time interval into N steps of size $\epsilon = T/N$, where N is large and will later be sent to infinity. That is actually already a lattice in time. This yields

$$\begin{aligned} U(a, b, T) &= \sum_j \sum_{i_j} \langle a, t_N | e^{-iH\epsilon} | i_{N-1}, t_{N-1} \rangle \dots \langle i_1, t_1 | e^{-iH\epsilon} | b, t_0 \rangle \\ &= \int \prod_i dq_i \langle q_a, t_N | e^{-iH\epsilon} | q_{N-1}, t_{N-1} \rangle \dots \langle q_1, t_1 | e^{-iH\epsilon} | q_b, t_0 \rangle, \end{aligned} \quad (2.2)$$

where in the second line the result was rewritten in terms of a set of continuous eigenstates of the (generalized) position operator Q_i . These are therefore $N - 1$ integrals.

If, as is the case for all systems of interest in the following, the Hamiltonian separates as

$$H = \frac{1}{2} P_i^2 + V(Q),$$

where the P_i and Q_i are the M canonically conjugated momenta, then for ϵ arbitrarily small the Baker-Campbell-Hausdorff formula

$$\exp F \exp G = \exp \left(F + G + \frac{1}{2}[F, G] + \frac{1}{12} ([[F, G], G] + [F, [F, G]]) + \dots \right).$$

yields

$$e^{-iH\epsilon} \approx e^{-\frac{i\epsilon}{2}P_i^2} e^{-i\epsilon V},$$

i. e. for infinitesimally small time steps the exponentials can be separated. Assuming the states to be eigenstates of the position operator and furthermore inserting a complete set of (also continuous) momentum eigenstates allows to rewrite the transition matrix elements as ordinary functions

$$\langle q_{i+1}, t_{i+1} | e^{-iH\epsilon} | q_i, t_i \rangle = e^{-\epsilon V(q_i)} \int \prod_j \frac{dp_j^i}{2\pi} \prod_k e^{-i\epsilon \left(\frac{p_k^i{}^2}{2} - ip_k \frac{q_k^{i+1} - q_k^i}{\epsilon} \right)}, \quad (2.3)$$

where products run over the number of independent coordinates M . This infinitesimal step is also known as the transfer matrix, which transfers the system from one time slice to another. In fact, even if the Hamilton operator is not known, but only the transfer matrix, it is possible to construct the full transition amplitude, as this only requires to sample all possible transfer matrices at every time slice. This possibility will be useful later on.

Defining

$$\mathcal{D}p\mathcal{D}q = \prod_i^N \prod_j^M \frac{dp_j^i dq_j^i}{2\pi}, \quad (2.4)$$

and thus in total $2NM$ integration measures yields the first formulation of the path integral

$$U(a, b, T) = \int \mathcal{D}p\mathcal{D}q e^{-\epsilon p_k \frac{q_j^{i+1} - q_j^i}{\epsilon}} e^{-i\epsilon H(p_j^i, q^j)}$$

Defining

$$\frac{q_j^{i+1} - q_j^i}{\epsilon} = d_t q_j^i$$

and performing the Gaussian integrals over the momenta yields

$$U(a, b, T) = \int \mathcal{D}q e^{i \sum^N \epsilon L(q_j^i, d_t q_j^i)} \stackrel{N \rightarrow \infty}{=} \int \mathcal{D}q e^{iS}$$

$$\mathcal{D}q = \prod_i^N \prod_j^M \frac{dq_j^i}{\sqrt{2\pi\epsilon}},$$

where L is the Lagrange function of the system, thus arriving at the original idea (2.1).

Considering the result in detail, it is important to note one important feature. The definition requires to chose any straight line between every point on every of the time

slices. Thus, in general paths will contribute which are not differentiable. This is a very important insight: Quantum physics differs from classical physics not only by including all possible paths, but also by including not only differentiable paths. This is in stark contrast to Hamilton's principle of classical mechanics.

Passing now to a field theory, the transition is the same as in classical mechanics: The paths are replaced by the fields, the Lagrange function by the Lagrangian density, and the action is an integral over space-time. Of particular importance is then the partition function

$$Z = \int \mathcal{D}\phi e^{i \int d^d x \mathcal{L}(\phi, \partial_\mu \phi)}, \quad (2.5)$$

where the integral is over all possible field configurations, i. e. all possible values of the field at all space-time points, including any non-differentiable ones². Since any field configuration includes the time-dependence, the path-integral can be considered as an integral over all possible histories of the universe described by the Lagrangian \mathcal{L} , from the infinite past to the infinite future. Thus, the path integral makes the absence of locality in quantum physics quite manifest. The partition function (2.5) is essentially the transition function from the vacuum to the vacuum. It is important to note that in the whole setup the field variables are no longer operators, like in canonical quantization, but ordinary functions. Nonetheless, the name 'operator' has stuck to such objects in the lattice literature, and will also be used here.

While the vacuum-to-vacuum transition is a very useful quantity, what is really important are the expectation values of fields, the correlation functions. These can be determined in a very similar way as before to be

$$\langle T(\phi(x_1) \dots \phi(x_n)) \rangle = \int \mathcal{D}\phi \phi(x_1) \dots \phi(x_n) e^{i \int d^d x \mathcal{L}(\phi, \partial_\mu \phi)}, \quad (2.6)$$

where there are two important remarks. One is that also the fields in the integration kernel are ordinary functions. The other is that, independent of the ordering of the fields in the path integral, this expression automatically creates the time-ordered correlation functions.

2.2 Euclidean space-time lattices

While the formulation (2.6) of correlation functions is useful, it is in practice quite complicated to use, beyond perturbation theory, due to the oscillatory nature of the exponential.

²In fact, it can be shown that those are the dominating ones in the continuum. Making sense out of this expression in the continuum is highly non-trivial and requires to pass from Riemann integrals to different definitions of integrals, but this is not the subject of this lecture.

This makes especially any numerical treatment complicated up to the point of practically impossible.

This problem can be circumvented by changing from Minkowski space-time to Euclidean space-time. This is formally done by a Wick rotation, i. e. by analytically continuing from t to it . The resulting expression for the correlation functions is then given by

$$\langle \phi(x_1) \dots \phi(x_n) \rangle = \int \mathcal{D}\phi \phi(x_1) \dots \phi(x_n) e^{-\int d^d x \mathcal{L}(\phi, \partial_\mu \phi)},$$

creating an exponential damping. Also, there is no time-ordering in Euclidean space-time. This damping makes numerical evaluations feasible, as will be discussed at length in chapter 4.

In Euclidean space-time it becomes then a well-defined statement to introduce a lattice: Space-time is reduced to a finite volume, which is further reduced to a finite number of discrete points. This set of points is called the lattice. Usually, the lattice is hyperrectangular³, of lattice extent N_μ in the different directions, and the points are equally spaced in each direction with a so-called lattice spacing a_μ . The total lattice volume is then $\prod_\mu (a_\mu N_\mu)$. If all N_μ and a_μ are equal, this is called a hypercubic lattice.

This finiteness is also the key to make numerical treatments in chapter 4 possible. There is one thing then necessary: Boundary conditions, like for every finite volume. They have to be fixed. However, eventually the focus will be on very large volumes, and therefore the boundary conditions will become irrelevant. Therefore, unless noted otherwise, they will be chosen to be periodic for convenience.

In Fourier space, a finite extent implies that there is a smallest absolute value of a non-zero momentum possible, given by $\pi/(a_\mu N_\mu)$ for any direction. The discreteness of the lattice will at the same time impose a maximal momentum of size π/a_μ for every direction. This is very similar to a crystallographic structure in solid state physics, and therefore also the name Brillouin zone is used for this momentum range. This Brillouin zone is given by $-\pi/a_\mu < p_\mu < \pi/a_\mu$, as $\exp(i2\pi x_\mu/a) = 1$. With periodic boundary conditions, this works indeed like in the solid state case. However, this also implies that results at positive and negative momenta are related by periodicity, and only half of the information is independent.

Thus, besides making the theory numerically tractable, the lattice introduces both an infrared cutoff as well as an ultraviolet cutoff in the theory. This makes all quantities well-defined on a lattice. Especially, in a quantum field theory all quantities are regularized by

³There have been many investigations of other geometries, including random lattices. So far, none of these have shown any advantage compared to a hyperrectangular lattice. Since a hyperrectangular lattice lends itself most easily to numerical aspects, this has become the standard discretization of space-time.

the lattice, and the lattice itself defines regularization.

The so-called thermodynamic limit is then defined as first sending the lattice spacing to zero and then the volume to infinity, where the order matters. This recovers the continuum theory. This will be discussed in more detail in section 2.7. The thermodynamic limit therefore corresponds to removing the regulator in a quantum field theory. Thus, before doing so, all quantities have to be renormalized, as otherwise they will either diverge or become zero. This points to a very important conceptual insight about lattice calculations. Because the number of lattice points is finite, a system has only a finite number of degrees of freedom. The theory at hand is therefore no longer a quantum field theory, but rather a quantum many-body system. In fact, a modern-day lattice calculation may involve often billions of degrees of freedom. Nonetheless, this number is still finite. Thus, lattice calculations are a possibility to approximate a quantum field theory by quantum mechanics. Since quantum mechanical systems are mathematically much better under control, it is possible to make much more powerful statements about systems on a lattice. In fact, it is even possible to prove that most interacting lattice theories are well-defined, which is not possible in the continuum. Also, many far-reaching exact statements can be made about lattice theories⁴.

There are two important remarks. One is that it appears at first worrisome that an analytical continuation is made while previously it was said that non-differentiable functions are involved. Fortunately, it can be proven, at least on a finite lattice, that this procedure is well-defined for any observable quantity, the so-called reconstruction theorem. However, the reconstruction of a single correlation function in Minkowski space-time usually requires an infinite number of Euclidean correlation functions, and thus a full reconstruction is in principle only approximately possible. As will be seen throughout this lecture, if only partial information is required, which is for many purposes usually more than sufficient, it is often possible to obtain a quasi-exact result.

The second is that when performing the analytical continuation the Lagrangian becomes the Hamiltonian, $T - V \rightarrow T + V$. Thus, the partition function becomes equivalent to the density operator of a statistical system in equilibrium. Hence, strictly speaking a computation in Euclidean space-time becomes a calculation in Minkowski space-time in

⁴This is one of the reasons that many people take lattice as a way to define quantum field theory, given that anyhow gravity will have to have a say about whether a thermodynamic limit is required. At the current time, this is a question of taste, as without a quantum theory of gravity there can be made no strong statement whether the actual thermodynamic limit is relevant for describing nature or not. Conceptually, of course, it would be very good to be able to perform it. However, as well as the structure of the theories are under mathematical control on a finite lattice, as little is the thermodynamic limit under mathematical control for essentially all interacting theories.

equilibrium at zero temperature and density in a $d + 1$ -dimensional system. In this case, the actual time becomes an auxiliary direction to implement the statistical nature of the operator, and the system is time-independent. The would-be time of the Euclidean lattice is actually a space-like, and thus Euclidean, direction of the $d + 1$ -dimensional system in Minkowski space-time. As the statement above shows, there is also another proof that all non-equilibrium correlation functions can be, with similar problems, be reconstructed from the equilibrium ones. Thus, no information is lost in this way, and this view is particularly useful to harvest the knowledge about numerical treatments of statistical systems in solid-state physics in chapter 4.

Hence, performing the Wick rotation is a well-defined, and very useful, trick.

2.3 Discrete analysis

On a finite lattice, analytic operations like integrating and taking derivatives become thrown back to their original definitions in terms of Riemann sums and finite differences. As a consequence, it may become important how these are implemented. As the theories live on a discrete lattice, and thus fields and other quantities are only known at fixed lattice points, any evaluation should only include the lattice points themselves.

As a consequence, derivatives can either be defined as a forward, backward or midpoint derivatives,

$$\partial_\mu \phi(x) \rightarrow \partial_\mu^f \phi(x) = \frac{\phi(x + e_\mu a) - \phi(x)}{a} \quad (2.7)$$

$$\partial_\mu \phi(x) \rightarrow \partial_\mu^b \phi(x) = \frac{\phi(x) - \phi(x - e_\mu a)}{a} \quad (2.8)$$

$$\partial_\mu \phi(x) \rightarrow \partial_\mu^m \phi(x) = \frac{1}{2a} (\phi(x + e_\mu) - \phi(x - e_\mu)), \quad (2.9)$$

respectively, where e_μ are the unit vectors in direction⁵ μ . While in the limit $a \rightarrow 0$ all are equivalent, they may introduce order $\mathcal{O}(a)$ discretization artifacts at any finite lattice spacings, which are different for the different versions. This problem no longer persist for the unique lattice Laplacian, which takes the form⁶

$$\partial^2 \phi(x) \rightarrow \partial_\mu^f \partial_\mu^b \phi(x) = \sum_{\pm\mu} \frac{\phi(x + e_\mu a) + \phi(x - e_\mu a) - 2\phi(x)}{a^2}. \quad (2.10)$$

⁵Note that because of the Wick rotation often μ is replaced by a Latin index and/or the directions are counted from 1 to 4 rather than 0 to 3. Neither of this will be done here.

⁶Summation on $\pm\mu$ implies a sum over all positive and negative directions, i. e. positive and negative values of the unit vectors. Using μ instead implies a sum only over positive directions of e_μ .

Likewise, an integral as a Riemann sum takes the form

$$\int d^4x \phi(x) \rightarrow a^4 \sum_x \phi(x),$$

falling back to its original definition. As a consequence, a functional integral becomes indeed again a product of normal integrals⁷,

$$\int \mathcal{D}\phi(x) \rightarrow \prod_x \int d\phi(x),$$

as x is just a counting variable in this context.

One apparent feature in all of these definitions is that the lattice spacing a is not only providing the finite differences, but also the dimension of the analytic operation. As a consequence, it is often useful to rescale all dimensionful quantities, e. g. fields or coupling constants, also by a to obtain dimensionless quantities. Then, all dimensions are reinstated by simply multiplying all quantities just by appropriate powers of a . This also implies that it is possible to measure all quantities in units of a , and thereby eliminating at fixed lattice spacings all explicit appearances of a by setting a to one. E. g., in this way

$$\partial^2 \phi(x) \rightarrow \sum_{\mu} (\phi(x + e_{\mu}) + \phi(x - e_{\mu}) - 2\phi(x)) = \sum_{\pm\mu} (\phi(x + e_{\mu}) + \phi(x - e_{\mu})) - 2d\phi(x),$$

where ϕ is now dimensionless and d is the dimensionality of space-time, usually 4 in the following. Of course, when taking the limit $a \rightarrow 0$ it is important to take particular care of it. However, since most of the following will happen on a finite lattice, and the continuum limit will only be occasionally appearing, a will henceforth be scaled out and set set to 1, except when noted otherwise.

2.4 The free scalar field

To exemplify this, consider the free scalar field. In the continuum, it is described by the Lagrangian

$$\mathcal{L} = \frac{1}{2} \partial_{\mu} \phi \partial_{\mu} \phi + \frac{M_0^2}{2} \phi^2.$$

Using the prescription of the previous section, the lattice version of this theory becomes

$$S = \sum_x \mathcal{L} = -\frac{1}{2} \sum_{x,\mu} \phi(x) \phi(x + e_{\mu}) + \frac{1}{2} (2d + M^2) \sum_x \phi(x)^2, \quad (2.11)$$

⁷Note that the continuum theory can, strictly spoken, not be taken as the limit of a product of Riemann integrals, but requires the more general concept of Ito integrals. This will be implicitly taken care of when including the exponentiated action as part of the measure.

where the periodicity of space-time has been used to rewrite the kinetic term, and $M_0 = aM$. The coordinate x is actually a vector of integers, which runs from 0 to $N_\mu - 1$ (or 1 to N_μ , depending on conventions).

This theory can be solved exactly. To do so, note that it is possible to map a vector of four integers x_i to a single integer, a master index n , e. g. by the prescription

$$n = x_0 + N_0x_1 + N_0N_1x_2 + N_0N_1N_2x_3, \quad (2.12)$$

of course, the choice of the fastest running index is arbitrary. The central point is that this permits to rewrite all expressions on a lattice as matrix-vector expressions⁸. Especially, the action (2.11) can be rewritten, now for $d = 4$, as

$$\begin{aligned} S &= \frac{1}{2}\phi_n K_{nm} \phi_m \\ K_{nm} &= -\sum_{\mu} (\delta_{n+e_{\mu},m} + \delta_{n-e_{\mu},m} - 2\delta_{nm}) + M^2\delta_{nm} \end{aligned} \quad (2.13)$$

where the notation $n \pm e_{\mu}$ means not addition by one, but a corresponding shift in the master index.

Since the action is quadratic in the fields, all possible N -point correlation functions can be evaluated exactly. All odd correlation functions vanish, and the even ones describe the propagation of $N/2$ non-interacting particles. Especially $N = 2$ is the propagator, which can be obtained according to the rules for multi-dimensional Gaussian integration as

$$\langle \phi_n \phi_m \rangle = K_{nm}^{-1}.$$

While the propagator K^{-1} could be obtained by direct inversion of K , it is useful to obtain it by Fourier transformation, using

$$K_{ln} K_{nm}^{-1} = \delta_{lm} = \int_{-\pi}^{\pi} \frac{d^d k}{(2\pi^d)} e^{ik(n-m)}.$$

The Fourier integral covers the whole momentum space of a finite lattice, the so-called Brillouin zone in analogy to solid state physics. All other momenta are determined due to periodicity. Moreover, periodicity implies that quantities are symmetric around zero momentum, and therefore the only independent (integer) momentum component run from 0 to $N_\mu/2$.

In Fourier space (2.13) becomes

$$K(k) = 4 \sum_{\mu} \sin^2 \frac{k_{\mu}}{2} + M^2,$$

⁸Even though this may numerically not always be the most efficient way.

where the dimensionless lattice momenta are given by

$$k_\mu = \frac{2\pi n_\mu}{N_\mu} \quad (2.14)$$

where the n_μ are the integer components. The first term of the denominator is therefore the lattice Laplacian in momentum space. There is no longer a second momentum argument, as momentum conservation applies⁹.

The inverse in Fourier space, and thus the lattice propagator of the free scalar field is thus

$$\langle \phi \phi \rangle(k) = D^t(k) = \frac{1}{4 \sum_\mu \sin^2 \frac{k_\mu}{2} + M^2}. \quad (2.15)$$

This explicitly demonstrates the periodicity of the lattice.

To recover the continuum limit requires to reinstantiate the dimensions. Comparison with the continuum expression implies that

$$K_\mu = \frac{2}{a} \sin \frac{(k_\mu/a)a}{2} \approx \frac{\pi n_\mu}{aN_\mu} + \mathcal{O}(a^2). \quad (2.16)$$

is the corresponding continuum momentum. k_μ/a is thus the lattice momentum. This already shows that lattice artifacts will deform the propagator if the approximation of the sine by its first Taylor term is not good.

To avoid these kind of lattice artifacts, it is possible to use quantities in which such artifacts are removed. Such quantities are called improved quantities. A trivial example is the dressing function of the propagator, defined as

$$Z(K^2) = \langle \phi \phi \rangle(k) \left(4 \sum_\mu \sin^2 \frac{k_\mu}{2} + M^2 \right) = 1,$$

instead of $\langle \phi \phi \rangle(k)(k_\mu^2 + M^2)$, which coincides with its continuum expression, because all lattice artifacts have been multiplied out. In this case, it is trivial. In the interacting theory, this is no longer obvious. Therefore, this would be considered a tree-level improvement, as it removes any artifacts at tree-level. Correspondingly, higher orders of lattice perturbation theory, to be discussed in section 2.6, can be used to do so beyond leading order. However, if the non-perturbative corrections to a quantity are large, there is no guarantee that a perturbative improvement will actually make the situation better, and may even make it worse at some fixed value of a . At any rate, it should be ensured that all improvements do not change the continuum result for the quantity in question, in this case the dressing function.

⁹As in solid state physics, up to translations of a whole Brillouin zone.

2.5 ϕ^4 theory on the lattice

As with continuum quantum field theory, ϕ^4 theory is also extremely useful to illustrate how to deal with interacting field theories on the lattice. For that purpose, its most simple implementation with a single, real degree of freedom will be used. As before in the free case, the lattice will be hypercubic with periodic boundary conditions in all dimensions.

In Euclidean space-time, the continuum action of this theory is

$$S = \int d^d x \left(\frac{1}{2} \partial_\mu \varphi \partial_\mu \varphi + \frac{m_0^2}{2} \varphi^2 + \frac{g}{4!} \varphi^4 \right). \quad (2.17)$$

Using the prescriptions of 2.3, the naive discretization of this theory is

$$S = a^4 \sum_x \left(\frac{1}{2} \partial_\mu^f \varphi \partial_\mu^f \varphi + \frac{m^2}{2} \varphi^2 + \frac{g^2}{4!} \varphi^4 \right), \quad (2.18)$$

and thus relatively straight-forward. Note that

$$am_0 = m, \quad (2.19)$$

i. e. the mass is rescaled such that the lattice mass m is dimensionless.

While this version is, indeed, correct, it turns out to be quite awkward to use. By performing

$$\varphi = \frac{2\sqrt{\kappa}}{a} \phi \quad (2.20)$$

$$a^2 m^2 = \frac{1 - 2\lambda}{\kappa} - 2d \quad (2.21)$$

$$g = \frac{6\lambda}{\kappa}, \quad (2.22)$$

the action takes the form

$$S = \sum_x \left(-2\kappa \sum_\mu \phi(x) \phi(x + \mu) + \phi^2 + \lambda(\phi^2 - 1)^2 \right), \quad (2.23)$$

where an irrelevant constant term has been dropped, and a partial integration has been performed in the kinetic term. The appearance of $2d$ in (2.21) stems from a rewriting of the discretized second derivative (2.10), where the local term has been combined with the mass term. The parameter κ is called the hopping parameter, in analogy to solid-state theories, where such terms described the probability of a quanta to move from one location to another (hop) in a (crystal) lattice. Its connection with the mass shows that the (inertial) mass opposes such hops. The rescaling also ensures that in the limit of $\lambda \rightarrow \infty$ the length

of the field is (classically) frozen to 1, irrespective of the other parameters of the theory¹⁰. Thus, the rescaling also rescales the amplitude. Note that nonetheless the amplitude of the field can range at every lattice point, at least in principle, within $(-\infty, \infty)$.

It is this theory for which now basic concepts on the lattice will be discussed. Though in the following many statements will be formulated for four dimensions, almost anything actually holds in arbitrary dimensions with minor modifications.

2.6 Lattice perturbation theory

2.6.1 Particle properties

While the true power of a lattice formulation is the possibility to perform non-perturbative calculations, it is also possible to do perturbation theory on a lattice. The reasons why this is nonetheless interesting are twofold.

One is that many non-perturbative calculations are performed using the numerical techniques of chapter 4. Due to the limited amount of computing time this often restricts the size of the lattice such that at most two orders of magnitude in scale differences can be covered, and sometimes much less. For theories which are asymptotically free in the infrared or ultraviolet it is possible to continue the results of the numerical calculation using perturbative means. This may even be useful for numerical calculations itself, as knowledge of perturbative renormalization properties can help to improve the calculations.

The second reason is that the lattice regularization provides a mathematically well-defined framework for quantum field theory. The renormalization, and even definition, of interacting quantum field theories in the continuum is far less well-defined in a mathematical sense. Thus, statements using perturbation theory on the lattice are often mathematically more reliable¹¹.

It is useful to have first a look at the free case. Though it can be solved exactly, it provides a lot of insight on how perturbation theory works on the lattice.

The starting point is the propagator (2.15). First consider the pole mass of the particle.

¹⁰In the particular case of a real field, its value is then restricted to ± 1 , and the theory becomes a spin system. That is, however, particular to the real case, and does not happen for more than one internal degree of freedom, especially the most relevant case of four real degrees of freedom.

¹¹However, the question whether they hold in the continuum limit is much less under mathematical control, and many statements are only true on a finite lattice. Given the aforementioned argument that in absence of quantum gravity this may be anyway an irrelevant argument, this will not be dwelt too much upon here.

It is located at the (analytically continued) poles of the propagator, which satisfy

$$\begin{aligned} P_4 &= \pm i(m^2 + \vec{P}^2)^{\frac{1}{2}} \\ P_\mu &= 2 \sin \frac{p_\mu}{2}, \end{aligned}$$

where the p_μ are again the dimensionless lattice momenta (2.14), and for the moment the formulation (2.18) of the action is again used. This implies that the lattice energy of the particle, defined as

$$p_4 = \pm i\omega$$

satisfies

$$\begin{aligned} 2 \sinh \frac{\omega}{2} &= \sqrt{m^2 + \vec{P}^2} \\ \cosh \omega &= 1 + \frac{1}{2}(m^2 + \vec{P}^2). \end{aligned}$$

The lattice energy-momentum relation therefore takes the form

$$\omega = \cosh^{-1} \left(1 + \frac{1}{2}(m^2 + \vec{P}^2) \right) = 2 \ln \left(\sqrt{1 + \frac{m^2 + \vec{P}^2}{4}} + \frac{1}{2} \sqrt{m^2 + \vec{P}^2} \right).$$

The dispersion relation on the lattice is therefore substantially different than in the continuum, where this would be just the ordinary Einstein formula $E = \sqrt{m^2 + \vec{P}^2}$. This will be important when analyzing systematic errors of calculations in section 4.5, and will be another example for the importance of lattice perturbation theory.

At small momentum this expands to

$$\omega = M + \frac{\vec{P}^2}{2m_*} + \mathcal{O}(\vec{p}^4),$$

and therefore the non-relativistic dispersion relation on the lattice is recovered. Note that this is equivalent, up to the rest mass, to the trivial dispersion law of a free particle in a crystal lattice in solid state physics. However, the two appearing masses are not just m , but have the values

$$M = 2 \ln \left(\sqrt{1 + \frac{m^2}{4}} + \frac{m}{2} \right) \approx m - \frac{m^3}{24} + \mathcal{O}(m^4) \quad (2.24)$$

$$m_* = \sinh M = m \sqrt{1 + \frac{m^2}{4}} \approx m + \frac{m^3}{8} + \mathcal{O}(m^4). \quad (2.25)$$

Thus the physical, i. e. rest mass M and the inertial mass m_* differ from the dimensionless tree-level mass m if the mass m is large. However, the lattice mass m is, due to (2.19),

exactly then not small if a is large, and thus the theory is far from the continuum limit. The relations (2.24) and (2.25) therefore describe how the finite lattice spacing affects the mass of the particle. They therefore form the lattice masses. Hence, the expansion should be considered rather as an expansion in a than an expansion in m .

This shows another use of lattice perturbation theory. With it, it is possible to estimate the influence of the non-zero lattice spacing, and also volume, on quantities measured on the lattice. Of course, analytical expressions will not be obtainable for genuine non-perturbative quantities, but in theories with asymptotic freedom this already helps a lot.

It is also interesting to compare the wave-function renormalization and the residuum at the pole. By construction, the wave-function renormalization is $Z_R = 1$. The residuum at the pole, and thus at the mass of the particle, is

$$Z_3 = \frac{M}{\sinh M} \approx 1 - \frac{m^2}{6} + \mathcal{O}(m^4)$$

and therefore approaches Z_R in the continuum limit, as it must be.

2.6.2 Interactions

To simplify the following, it will be assumed that no spontaneous symmetry breaking occurs¹².

In principle, the procedure works in a very similar way as in continuum quantum-field theory, i. e. the derivatives of the Lagrangian are used to define the Feynman rules. The differences arise from the fact that now the discretized lattice Lagrangian (2.18) is used, rather than the continuum version (2.17). While the differences will be found to be similar to the free case for the scalar field, the situation becomes more involved when non-scalar fields are involved, as will be discussed in section 5.6 for Yang-Mills theory.

The actual differences in Feynman rules are

- The tree-level propagator is (2.15), rather than the continuum one
- Momentum conservation for the lattice momenta (2.14) holds only up to multiples of 2π , just as in a crystal, due to the periodic boundary conditions
- Loop integrals are performed over the first Brillouin zone only, and are sums over the discrete lattice momenta (2.14). Note that momentum sums inside an integral can

¹²Which, in a full non-perturbative simulation, actually requires additional measures to have anyways, as discussed in section 4.8.

then be wrapped again into the first Brillouin zone by periodicity. This automatically regulates all loop integrations both in the infrared and the ultraviolet

- The rescaling of the fields (2.20) requires a rescaling of any n -leg diagram by $(2\kappa)^{-\frac{n}{2}}$

The separation in 1PI and amputated diagrams is done as in the continuum, and the symmetry factors are the same. Also the vertex is the same as in the continuum, and given by $-g$.

As an example, consider the (resummed) propagator D to two loops. This involves three diagrams, and thus three contributions. These are the tadpole diagram, the stacked tadpole diagram, and the sunset diagram,

$$\begin{aligned} -\frac{1}{2\kappa D}(p) &= -(P^2 + m^2) - \frac{g}{2}J_1 + \frac{g^2}{4}J_1J_2 + \frac{g^2}{6}I \\ J_n &= \sum_q (D^{\text{tl}})^n \\ I &= \sum_q \sum_k D^{\text{tl}}(q)D^{\text{tl}}(k)D^{\text{tl}}(p - q - k). \end{aligned}$$

Now, consider the situation to leading order in g . Then only the leading tadpole J_1 contributes. As it is momentum independent, it is only a shift in the mass. Requiring

$$D = \frac{1}{P^2 + m_R^2}$$

to define the renormalized mass to this order leads to an expression for the pole mass M^g to order g very much like (2.24),

$$M^g = 2 \ln \left(\sqrt{1 + \frac{m_R^2}{4} + \frac{m_R}{2}} \right) = M^0 + \frac{1}{2m_*} \frac{g}{2} J_1, \quad (2.26)$$

where M^0 is the tree-level mass (2.24) and m^* the inertial tree-level mass (2.25).

This leaves to calculate the contribution J_1 . For simplicity, assume that the lattice size¹³ $N \rightarrow \infty$, and thus the summation is actually an integral over the first Brillouin zone,

$$J_1 = \int_{-\pi}^{\pi} \frac{d^4q}{(2\pi)^4} \frac{1}{4 \sum_{\mu} \sin^2 \frac{k_{\mu}}{2} + m^2} \stackrel{m=0}{=} r_0 \approx 0.155. \quad (2.27)$$

The limit $m \rightarrow 0$ is actually interesting as, as will be discussed in the next section 2.7, this is actually necessary close to the continuum limit. Expanding the integrand in powers of

¹³This is often not a good approximation, but without it it is usually not possible to obtain analytical results at all.

m^2 yields in the next order in m

$$\begin{aligned} J_1 &= r_0 + m^2 \left(\frac{1}{16\pi^2} \ln m^2 + r_1 + \mathcal{O}(m^2) \right) \\ r_1 &\approx -0.0303, \end{aligned} \quad (2.28)$$

and thus the integral is actually well-behaved for $m \rightarrow 0$.

On the lattice, the mass renormalization can also be viewed from a different perspective. The renormalized mass, with a reinstated, reads

$$(am_R)^2 = (am)^2 + \frac{gr_0}{2} + \frac{g}{32\pi^2} (am)^2 \ln(a^2 m^2) + \frac{gr_1(am)^2}{2} + \mathcal{O}(g^2, (am)^4). \quad (2.29)$$

To obtain a finite, renormalized, mass in the continuum limit $a \rightarrow 0$, this implies that it is necessary to tune the bare parameter m accordingly to¹⁴

$$(am)^2 = -\frac{g}{2} r_0 + \mathcal{O}(g^2),$$

and then the renormalized, physical mass becomes finite, and in this case actually zero. If there should remain a non-zero value, it is necessary to include it in this condition. Thus, to obtain a definite mass in the continuum limit requires to tune the bare mass accordingly. This is, however, nothing else but the usual tuning of counter terms in continuum calculations.

It should be noted that (2.27) appears to be never divergent, as (2.29) is always finite, even for $a \rightarrow 0$. The reason is that the actual limits in (2.27) are always finite in units of a . In fact, reinstating units in (2.27) yields upper and lower limits of $\pm\pi/a$. Thus, the expression (2.29) is not only regulated, but by expressing everything in units of a already renormalized. The divergences would pop up when introducing a again in (2.28), which would give a leading term r_0/a^2 , which diverges in the continuum limit.

As in lattice calculations mostly the form (2.23) rather than (2.18) is used, this translates into a tuning of the hopping parameter (2.21)

$$\kappa = \kappa_c = \frac{1}{8} + \left(3r_0 - \frac{1}{4} \right) \lambda + \mathcal{O}(\lambda^2). \quad (2.30)$$

This implies that by choosing $\kappa = \kappa_c$, the critical hopping parameter, a particular mass, in this case zero, is fixed for the particle in the continuum limit. This is, unfortunately only half the truth. This would be correct, if all relevant contributions would be entirely due to the leading perturbative contribution, and the rest could be neglected. But this is actually

¹⁴Note that by this tuning all other terms become automatically of higher order in g and can therefore be neglected.

a fine-tuning problem. If not the contributions to all orders, and also all non-perturbative contributions, are exactly canceled, any surviving term will act like a non-zero mass. Since therefore the actual value cannot be determined exactly, as otherwise the theory would have already been solved and lattice simulations would be pointless, this implies that in actual numerical calculations the correct value for the bare parameters need to be determined in an iterative procedure: Choose a value for κ , perform (numerically by extrapolation) the limit $a \rightarrow 0$ for the renormalized mass, improve the choice of κ and repeat, until the desired mass has been hit satisfactorily good.

It is interesting to briefly discuss the opposite limit, i. e. $\lambda \rightarrow \infty$. In this case the potential in (2.23) enforces that the length of ϕ is frozen to 1, and thus the only possible values are $\phi = \pm 1$. Other than that the term does not contribute. The action then reduces to

$$S = -2\kappa \sum_x \sum_\mu \phi(x)\phi(x + \mu).$$

This is then just the ordinary Ising model. Thus, the infinite-coupling limit of the ϕ^4 theory is reducing to a quantum-mechanical system. It should be noted that this limit does not easily commute with the continuum limit, and care has to be taken, and this statement is therefore on a finite lattice.

2.6.3 Hopping expansion

The formulation (2.23) allows also for a quite different expansion than in λ : An expansion¹⁵ in κ . As κ is a measure of how easy it is to move from one site to another, this is also called a hopping expansion. Assume therefore that κ is small, but not necessarily that λ is small. Because of (2.21) this implies that the mass is large, at least compared to the coupling λ . Note, however, that this refers to the lattice mass. As discussed in section 2.7, the continuum limit requires a vanishing lattice mass. Therefore, the hopping expansion is automatically dealing with an extremely coarse lattice. While it can therefore not be expected that this will give results quantitatively agreeing with the continuum it is possible that qualitative results may still agree, under conditions formulated in section 2.7.

In the extreme case of $\kappa = 0$, (2.23) reduces to

$$S = \sum_x (\phi^2 + \lambda(\phi^2 - 1)^2) = \sum_x S(x).$$

In this case, the path-integral can be calculated exactly, since the Boltzmann weight

¹⁵This is sometimes also called a high-temperature expansion, as formally κ plays the same role as the inverse temperature of a statistical system, see chapter 7.

factorizes,

$$Z = \int \mathcal{D}\phi e^{-S} = \prod_x \int \mathcal{D}\phi e^{-S(x)} = \left(\int \mathcal{D}\phi e^{-S(x)} \right)^V = Z_I^V,$$

where V is the lattice volume. While still not exactly solvable, the single-site integral can be determined numerically to essentially arbitrary precision.

Corrections beyond κ can then be determined in a series expansion in κ . To see how this works, rewrite the term proportional to κ in (2.23) as

$$S_I = -2\kappa \sum_{nn} \phi(x)\phi(y),$$

where nn now indicates that the sum should be performed over nearest neighbors, and that x and y only differ by a single lattice spacing. The partition function then reads

$$Z = \int \mathcal{D}\phi \prod_z e^{-S(z)} \prod_{nn} e^{2\kappa\phi(x)\phi(y)}, \quad (2.31)$$

where the site-local and nearest-neighbor terms have been separated.

Expand now the second factor in κ ,

$$\prod_{nn} e^{2\kappa\phi(x)\phi(y)} = \prod_{nn} \left(1 + 2\kappa\phi(x)\phi(y) + \frac{(2\kappa)^2}{2!} (\phi(x)\phi(y))^2 + \mathcal{O}(\kappa^3) \right).$$

As is actually quite often the case, it is now possible to give this object an interpretation in graph theory. For this, note that $\phi(x)\phi(y)$ can be considered to be a (directed) connection from site x to site y . Performing the product then creates products of such connections, creating a bond. All such bonds start at some initial point i and end at some final point f . If there is more than one bond between two points, their multiplicity is $m(i, f)$. A graph \mathcal{G} is then the collection of all bonds with the same initial point and final point. The number of bonds in a graph is called by $L(\mathcal{G})$. The number of points ending in a lattice point, called in this context a vertex ν , is called $N(\nu)$.

It can then be shown that

$$\begin{aligned} \prod_{nn} e^{2\kappa\phi(x)\phi(y)} &= \sum_{\mathcal{G}} (2\kappa)^{L(\mathcal{G})} c(\mathcal{G}) \prod_b [\phi(i(b))\phi(f(b))] \\ c(\mathcal{G}) &= \prod_{if} \frac{1}{m(i, f)!}, \end{aligned}$$

and thus an expression determined entirely in terms of possible graphs, which can be created on the lattice. That this is possible is actually not specific to ϕ^4 theory, but rather generic: Lattice actions can be rewritten in geometric terms describing geometric objects on the lattice. This is also true for continuum perturbation theory, though in that case as

graphs the Feynman graphs serve. The underlying structure in both cases is exactly the same: The fact that they are a combinatorial rewriting of expressions of a infinite product of an infinite sum.

This implies that to evaluate (2.31), it is necessary to evaluate integrals of type

$$\gamma_k = \frac{1}{Z_1} \int \mathcal{D}\phi \phi^k e^{-S_I}.$$

Because S_I is an even function $\gamma_{2n+1} = 0$. Note that this implies that only closed bonds will contribute. Open bonds have at their endpoints fields with an argument not reappearing, and therefore they vanish. Closed loops have the endpoints squared, and will therefore contribute. For even n , this can still not exactly be integrated, but calculated numerically to essentially arbitrary precision.

This allows to finally rewrite (2.31) as

$$\frac{Z}{Z_1^V} = \sum_{\mathcal{G}} (2\kappa)^{L(\mathcal{G})} c(\mathcal{G}) \Pi_{\nu} \gamma_{N(\nu)}.$$

This is an exact expression, as long as the lattice spacing remains finite. In the limit $a \rightarrow 0$, the sums become integrals, and the exchange of summation and integration may no longer be possible. However, it appears that this is often a surprisingly good approximation to the continuum theory.

To give an example of how a result could look like, the normalized partition function up to order κ^4 in d dimensions is

$$\frac{Z}{Z_1^V} = 1 + (2\kappa)^2 \frac{B d \gamma_2^2}{2} + \mathcal{O}(\kappa^4),$$

with B a constant. While this result is particular for the partition function, it is possible to derive similar results for other quantities, like the free energy¹⁶ or for 1PI correlation functions, and essentially any quantity of interest.

2.7 The phase diagram and the continuum limit

After having introduced the first non-trivial examples of a theory on the lattice, the next logical question is how to obtain the continuum quantum field theory. It is a long discussion whether this is actually necessary, as any theory requiring normalization anyways will break down at some scale Λ , and it would therefore be sufficient to consider the situation

¹⁶In this case the expansion is called linked-cluster expansion for reasons stemming from solid-state physics.

at finite lattice spacing. However, only for theories for which the continuum limit can be taken, at least in principle, it can be ensured that the result is independent of the regulator, and therefore whether results are obtained using a lattice or, say, some Pauli-Villars regulator.

The necessary procedure can already be discussed when considering the free scalar field of section 2.4. The single parameter of the theory was the mass, which had the property that the continuum mass is given by the lattice mass by

$$M_0 = \frac{M}{a}.$$

This implies that having a to zero is achieved by sending the lattice parameter M to zero, since

$$a = \frac{M}{M_0}$$

and the continuum, dimensionful mass M_0 is fixed. However, for the free theory on the lattice $1/M$ is just the correlation length. Thus, the continuum limit requires the correlation length to go to infinity, and thus corresponds to a second-order phase transition of the discrete lattice model. Of course, on a finite volume there is no true phase transition, and thus the order of the two limits is important here.

At any rate, the important insight is that the continuum limit of a lattice field theory is defined by a second-order phase transition in the phase diagram spanned by the parameters of the theory. Thus, the existence of a second-order phase transition is necessary for a lattice field theory to have a continuum limit. If the phase diagram of a lattice field theory does not contain a second-order phase transition, it is not possible for the theory to have a continuum limit. This corresponds in the continuum limit to the impossibility of removing the regulator.

In such a case, there exists no continuum field theory. However, this does not imply that the theory makes no sense at all. If in some regions of the phase diagram the dependence on the lattice spacing is small, the theory may still serve well as a low-energy effective theory, as is e. g. also possible for some non-renormalizable continuum field theories, as long as the dependence on the regulator is sufficiently small for the purpose at hand.

When the theory becomes interacting, there are two further subtleties to be considered. One is that the parameters of the Lagrangian are often subject to renormalization, and thus would go to either zero or infinity when sending the regulator to infinity, which corresponds to sending the lattice spacing to zero. If this is the case, they are not useful to identify a correlation length, as their size depends on the renormalization scheme. Thus, it is necessary to consider renormalization scheme and scale invariant quantities, e. g. the proton mass. With it, in the same way as above a correlation length and a second order

phase transition can be identified. Thus, even for interacting theories it is necessary and sufficient to search for a second-order phase transition in the quantum phase diagram of the theory.

The second subtlety is that a theory which interacts at a finite lattice spacing may cease to interact at a second-order phase transition, i. e. all correlation functions take the form of a non-interacting theory. Such theories are termed trivial, and e. g. the ϕ^4 theory of section 2.5 belongs very likely to this class, as well as QED.

Though the prescription to find the continuum limit for any given lattice field theory is thus straightforward, it is in practice far less straightforward, especially if the quantum phase diagram has a significant number of dimensions, where significant can already be two, depending on the theory: There is no analytic way to determine for any reasonable complicated, and for particle physics relevant, theory the phase diagram analytically¹⁷. It is therefore necessary to search for them using numerical methods, as discussed in chapter 4. As any numerical statement, no identification of a possible phase transition is therefore ever proof, and thus the (non-)existence of second order, and thus continuum, field theories stays therefore in most cases an assumption backed by (circumstantial) numerical evidence.

However, luckily, most interesting theories belong to the class mentioned above where the dependence on the lattice spacing, and thus the distance to the continuum becomes quickly practically negligible, and therefore no matter the status of the theory as a continuum field theory for most practical purposes the lattice theory at finite a is completely sufficient. As the general attitude in field theory is anyhow based on a hierarchy of effective field theories, until a better concept emerges, the problem of actually proving the existence of a second-order phase transition is usually not too relevant. However, it remains still a numerical challenge, which rarely can be augmented by analytical statements, to find regions of the phase diagram where a is sufficiently small to make practical applications possible.

There is, however, an important secondary problem. While the prescription above determines how to find a continuum limit, it does say nothing about the parameters of the theory for which the continuum limit has been reached. As is visible, e. g., in expressions like (2.30), the bare parameters of the Lagrangian need to be tuned very precisely to obtain a particular value for the renormalized parameters of the theory in the continuum limit.

To control the continuum limit of which theory is taken, it is therefore necessary to ensure that the limit is taken at fixed theory. For this purpose, consider a sequence of values for the lattice spacing a_n , with $a_{n \rightarrow \infty} = 0$. To actually deal with the desired theory in the continuum limit requires to keep all independent parameters the same for every

¹⁷Though some quite interesting statements can be made in general.

a_n . In a renormalizable theory, these are a finite number of parameters, e. g. for the ϕ^4 theory the renormalized mass and the renormalized interaction strength at every a_n . Note that the bare parameters κ and λ change and thus $a_n(\kappa_n, \lambda_n)$. This sequence defines a so-called line of constant physics, or sometimes also renormalization trajectory. Thus, the prescription for the continuum limit to a predetermined theory is to follow a line of constant physics to the continuum limit, rather than just perform the continuum limit.

While this seems simple enough, especially with expressions like (2.30) at hand, this is actually not so. First, as already noted below (2.30), we do not know the exact expressions to have at fixed a the desired physical properties, as otherwise we would have already solved the problem. Thus, in general, this requires an iterative procedure to find a line of constant physics. The second problem is that not for all possible desired values there is a second-order phase transition, and thus not all lines of constant physics end at such a transition. This may even differ within a theory, e. g. for the ϕ^4 theory there is overwhelming evidence that only the trajectory $g = 0$ ends in a second-order phase transition, and the theory is thus trivial.

This also implies that the line of constant physics can lead to very different points in the phase diagram, and structures like isolated points of second-order phase transitions, but also hypersurfaces are possible. Thus, the structure may indeed be very rich. And complicated to determine.

There is an important problem when following a line of constant physics. The lattice phase diagram, cannot only have such second-order transitions, but also phase transitions of other order, e. g. first order. This corresponds to lattice transitions, akin to quantum phase transitions of crystals in solid state physics. Since across such transitions quantities change non-analytically, so do all quantities measured along such a line of constant physics. Hence, it is, mathematically, not possible to continuously interpolate, or extrapolate, across such a non-analyticity. An extrapolation towards the continuum limit is then, in general, impossible. Thus, to extrapolate from results at a finite lattice spacing to zero lattice spacing it is required that no non-analyticities along the line of constant physics are encountered.

The serious problems for an actual calculation is, however, that for any non-exactly solvable model the presence or absence of such non-analyticities can not be detected without actually finding them. Therefore, for any interesting theories, the absence of non-analyticities along the lines of constant physics must actually be assumed, making the continuum limit generally not systematically controllable. That should always be kept in mind.

This implies that results from the hopping expansion of section 2.6.3 can usually not

be expected to connect to the continuum limit¹⁸.

Still, this may in actuality not be as bad as it seems. Quite often non-analyticities seem to happen only while a is still quite large and thus the lattice coarse. Also, often the consequence of non-analyticities seem to be not be a strong quantitative influence on many observables of interest. But it still requires care.

2.8 Detecting phase transitions

It has now be very often allured to phase transitions. However, as known from thermodynamics, there are no true phase transitions, i. e. no non-analyticities, in a finite volume, and thus on a finite lattice. Since practically, especially numerical, calculations are usually restricted to a finite volume the question arises how to then detect phase transitions. The answer to this lies in the way how phase transitions evolve when the volume is send to infinity.

Before doing so, an important distinction of two cases has to be made: Phase transitions which have or have not a finite correlation length, i. e. first or second order phase transitions. As has been discussed in section 2.7, only second-order phase transitions are really interesting, as they are associated with the continuum limit. However, first order are not, but can still exist, and make an extrapolation to infinity impossible. If a (first-order) phase transition happens at fixed (lattice) parameters, it is also called a bulk phase transition.

The tool to identify phase transitions is by the scaling of relevant quantities with the volume, a so-called finite-size-scaling analysis. Phase transitions are signaled by non-analyticities in thermodynamic bulk quantities, like the free energy, heat capacities etc. in form of either jumps or singularities. On a finite volume, these non-analyticities are absent.

Consider the specific heat C , or any other quantity showing the same qualitative behavior at a phase transition. Using similar arguments as when obtaining critical exponents in temperature¹⁹, it can be shown that such quantities behave as a function of lattice extension with a critical behavior

$$C = L^\alpha C_0(L^\beta) \tag{2.32}$$

¹⁸In fact, there are more fundamental reasons speaking against this, as e. g., any interacting continuum theory should be non-analytic in its coupling, and therefore any expansion can never be exact. This is known as Haag's theorem.

¹⁹In fact, as discussed in chapter 7, temperature is actually nothing but a finite time extent.

with $C_0(0)$ constant and $\alpha > 0$ and $\beta < 0$. The critical exponents α and β are characteristic for the universality class of the second order phase transition at hand. Like for the conventional critical exponents of the temperature behavior, (hyper)scaling relations between different exponents of different quantities exist.

Thus, identifying a second-order phase transition can be done by measuring a quantity like (2.32) as a function of L , and fitting the exponent, and extrapolating it to infinite volume. There is, however, one important issue to take into account. The relation (2.32) is true around the infinite-volume divergence. At $L = \infty$, this peak is located in the phase diagram at a critical value of the lattice parameters, say κ_c and λ_c for the ϕ^4 theory. However, on a finite lattice, the maximum of (2.32) may not be at the same values of κ and λ , as there are additional finite contributions,

$$C = L^{\alpha(\kappa_c, \lambda_c)} C_0(L^\beta, \kappa_c, \lambda_c) + C_1(\kappa, \lambda, L),$$

where C_1 is finite even for $L \rightarrow \infty$. The dependence of κ_c and λ_c have been added, as in principle there can be multiple phase transitions in the system at different positions in the phase diagram. Thus, in general the peak will move, and this has to be taken into account when fitting.

In fact, also the critical values show a scaling behavior of type

$$\alpha_{\max}(L) = \alpha_c + \frac{a}{L^\gamma} + \alpha_1(L)$$

where α can be either κ or λ (or the parameters of a given model), a is a constant, $\gamma > 0$, and α_1 vanishes quicker than $1/L^\gamma$ for $L \rightarrow \infty$. In fact, $\gamma = \beta$ in many cases due to scaling relations.

Of course, it can be argued that because of (2.29) it would also be sufficient to just track the correlation lengths directly. This is true. Unfortunately, in practical calculations this is usually much harder than monitoring bulk quantities. The scaling with volume of extensive quantities makes them much less sensitive to fluctuations than correlations.

Another possibility appears to be order parameters, if there is one associated with the second-order phase transition²⁰. However, as will be discussed in more detail in sections 2.9 and 4.8, this is not directly possible.

First order transitions are not associated with singularities, but with jumps, due to the existence of a latent heat. Also these jumps get washed out on a finite lattice. Similar as to the case of a second order phase transition, they can be identified by tracking some suitable quantity with a jump. A reasonable estimate of its behavior is, e. g.,

$$C(L, \alpha) = C_0 \tanh(L^d(\alpha - \alpha_c))$$

²⁰Second-order phase transitions without associated symmetry and order parameter are rare, but do exist.

where d is the dimension. This scaling with L^d is typical for a first order transition, though the tanh-behavior is not necessary. Note also that the jump position is in this case not changing with L , a typical signal of a first-order bulk transition. Thus, measuring a dependence on L^d in this way provides a possible tag for a first-order transition.

There are other, partly even more powerful, possibilities to identify a first-order transition by the fact that a coexistence of both phases and a hysteresis is possible. This will be discussed in more detail in section 4.7.

Note that there are also phase transitions with finite correlation length, but singularities only in even higher derivatives than susceptibilities²¹, including those which have such a singularity at an infinite number of derivatives. Since these are nonetheless signals of a non-analyticity in the partition sum, these also make any analytical continuation impossible. However, they are much harder to detect, as statistical fluctuations increase with the order of derivative, though in principle a higher-order finite-size-scaling analysis is also possible in this case. A particular example of such a case is the so-called roughening transition of the non-Abelian gauge theories, to be discussed in chapter 5, which has a singularity only at infinite order, but nonetheless effectively blocks an extrapolation from the equivalent of the hopping expansion in such theories.

2.9 Internal and external symmetries on the lattice

On the lattice, symmetries work somewhat different than in the continuum, especially external symmetries. But there are also subtle differences for internal symmetries, due to finite volume.

2.9.1 External symmetries

External symmetries are the Poincare symmetry, as well as symmetries connected to them, like chiral symmetry and supersymmetry. Chiral symmetry will be discussed in chapter 6, but supersymmetry is beyond the scope of this lecture. Here, Poincare symmetry will be the central issue.

There are two central points to be considered when it comes to Poincare symmetry. The fact that the lattice formulation is in Euclidean space-time, and on a lattice, see section 2.2.

²¹In the old Ehrenfest classification, these are phase transitions of higher-than-second order.

2.9.1.1 Translation symmetry

On a finite lattice translation symmetry is no longer a continuous symmetry, and it is only possible to translate by one unit of lattice spacing a . The translation group therefore becomes a discrete group. If the theory is endowed with boundary conditions, this group becomes finite, as, if t is the translation operator and there are N lattice sites and the boundary conditions are periodic, $t^N = 1$ holds. Thus, the translation group becomes the group Z_N . Correspondingly, all representations have to fall into such representations. Especially, this implies that energy levels are necessarily discrete, and there exists, as in solid state physics, the already noted Brillouin zone. Fourier transformation therefore maps to the dual space, like in a crystal.

2.9.1.2 Reflection positivity

Lattice theories are usually defined in Euclidean space-time. In section 2.2, it was stated that this is, generically, not a problem. While the arguments given there are adequate for standard quantum field theories, this is not generically true. As will be discussed in section 5.9, there are reasons to consider actions deviating from the natural ones. To ensure that they provide well-defined quantum-field theories is less trivial. The important criterion for this is that the transfer matrix (2.3) fulfills reflection positivity.

The reflection operator R on a finite lattice for a scalar field is defined as

$$\begin{aligned} R\phi(x) &= (\phi(Rx))^* \\ R(x_0, \vec{x}) &= (-x_0, \vec{x}), \end{aligned}$$

i. e. it is the equivalent of time reversal T in the continuum theory, and inherits from it its anti-linearity and product rule.

Consider now an arbitrary function

$$F = \sum_i \int dx_1 \dots dx_i f_i(x_1, \dots, x_i) \phi(x_1) \dots \phi(x_i) \quad (2.33)$$

where the functions f_i fall off sufficiently fast at infinity to make the integrals well-defined and have support only at positive times $x_0 \geq 0$.

Reflection positivity is now the statement than for any set of f_i

$$\langle (RF)F \rangle \geq 0 \quad (2.34)$$

holds. Defining the Schwinger function

$$\mathcal{S}_n(x_i, y_j) = \langle (\phi(y_1) \dots \phi(y_n))^\dagger \phi(x_1) \dots \phi(x_n) \rangle,$$

this amounts to

$$\sum_{ij} \int dx_1 \dots dx_n dy_1 \dots dy_n f(y_k)^* f(x_l) \mathcal{S}(Ry_k, x_l) \geq 0.$$

It thus means that any expectation value in the forward direction in time, smeared by a suitable test function, has a positive norm. It is thus essentially the same statement which is necessary to create a probability interpretation in Minkowski space-time.

On a finite, hypercubic²² lattice, there are now actually two possibilities to perform a reflection. It is possible to reflect in the plane $x_0 = 0$ and thus $x_0 \rightarrow -x_0$. This is called site reflection. Or in the plane between two lattice planes, e. g. at $x_0 = 1/2$, leading to $x_0 \rightarrow 1 - x_0$. This is called link reflection, as it happens in the region linking two lattice surfaces. Only if for both types of reflections (2.34) holds, reflection positivity holds, and the theory is well-defined.

Interestingly, and actually practically useful, is that if a theory is known to satisfy reflection positivity, it is possible to derive the Hamilton operator, and thus action, from the transfer matrix. Consider first site reflection, and a function F as defined in (2.33). This implies that a state created by F , $|\Psi\rangle = F|0\rangle$ satisfies

$$\langle \Psi | \Psi \rangle \geq 0,$$

nothing but a positive norm²³. This is true for any F . In particular, this implies it is also true, if F vanishes for a number n of lattice sites in positive time direction. Then, there is a gap between RF and F of $2n$ lattice sites, but their norm is still positive. A special case of such a function is any function F , which is translated by the transfer matrix T , (2.3), by n lattice spacings. This implies

$$\langle \Psi | T^{2n} | \Psi \rangle \geq 0.$$

This implies that T^2 is a positive operator, and that a Hamiltonian for a lattice translation by two units of lattice spacing can be defined as $T^2 = e^{-2aH}$, which is thus necessarily self-adjoint such that its exponential can be positive. In the same manner link reflections can be used to show that the same results holds if the initial separation of F is only one, but not two lattice spacings, due to the half lattice spacing displacement to start with. Thus, reflection positivity guarantees a positive norm of any set of states and the existence of an self-adjoint Hamilton operator.

It is interesting to apply link reflections to the scalar field to see how reflection positivity requires the hopping parameter to be positive. The only part of the action which has no

²²Similar arguments pertain to other structures

²³This already forebodes that gauge theories with their indefinite metric space are harder to deal with.

trivial properties under reflection is the kinetic energy in (2.23). Also here, almost all terms are not close to the reflection surface. For simplicity, drop the potential energy, as it will have trivial properties in the following. Then the action consisting only of terms involving only $x_0 \geq 1$ will be called S_+ . Likewise S_- collects all terms at negative times, and S_c the rest. It is then possible to decompose the action as

$$\begin{aligned} S &= S_+ + RS_+ + S_c \\ S_c &= -2\kappa \sum_x \phi(1, \vec{x})\phi(0, \vec{x}). \end{aligned}$$

Because $\phi(0) = R\phi(1)$, the Boltzmann weight of S_c can be expanded as

$$e^{-S_c} = \sum_n \kappa^n \sum_i c_{ni} \phi(1, i) R\phi(1, i),$$

where the c_{ni} are positive coefficients, and the i -notation indicates that the corresponding lattice sum needs to be evaluated. Using this, the relevant expectation value becomes

$$\langle FRF \rangle = \sum_n \kappa^n \sum_i c_{ni} \left| \int \Pi_{x_0 > 0} d\phi F\phi(1, i) e^{-S_+} \right|^2$$

where it has been used that the involved functions are real and known functions of the fields, and thus the action of R can be explicitly evaluated. This has to be positive for any operator F . This can only be true, if $\kappa \geq 0$. Thus, the hopping parameter needs to be positive, as announced.

2.9.1.3 Rotation symmetry and angular momentum

The lattice breaks the continuous (Euclidean) rotation symmetry $SO(d)$. For a hypercubic lattice, the remainder symmetry is the so-called hypercubic symmetry with symmetry group $H(d)$, a discrete group, known from crystallography.

As an example, consider two dimensions, yielding a square lattice with group $H(2)$. The only rotations possible are by $\pi/2$, or by multiples of this. Of course, there are also the parity transformations, i. e. P and T , which in two dimensions are also rotations. The group is actually isomorphic to Z_4 , as for any rotation $R(n\pi/2)^4 = 1$ is true, including P and T .

The most important consequence is that a discrete group has only a finite number of irreducible representations. This has extreme consequences when it comes to spin or angular momentum. Every spin value in the continuum, i. e. $s = 0, 1/2, 1, \dots$ is a separate irreducible representation of the rotation group.

As a consequence an infinite number of irreducible representations in the continuum are mapped to the same irreducible representation on the lattice. This implies that various spin values cannot be distinguished on the lattice.

Consider the most important case of four dimensions. The group $H(4)$ has 24 elements and 5 irreducible representations of dimensionality 1, 1, 2, 3, and 3, called A_1 , A_2 , E , T_1 and T_2 respectively.

Each of the irreducible representations is linked to a number of continuum spins. These are the sequence 0, 4,... for A_1 , 3, 6,... for A_2 , 2, 4, 5,... for E , 1, 3, 4, 5,... for T_1 and 2, 3, 4, 5,... for T_2 . Except for spins 0 and 1 all spin values appear in different representations. This implies that various components of the continuum spins are mapped to different discrete representations. Therefore, to collect all spin components for any spin larger than 1 will require to investigate more than one of the discrete representations. This is important, as the simplest possibility to identify what continuum spin is associated with a given object is to count degeneracies. This also implies that the, in the continuum limit different, parts of a given representation map to different subspaces of the continuum representations, but the mapping is not unique, as can be seen, e. g., for spin 4.

This can also be counted differently. The single component of spin 0 is located entirely in A_1 . All three components of spin 1 are located in T_1 . For spin 2, its five components are distributed over E and T_1 . For spin 3, its 7 components are distributed over A_2 , T_1 , and T_2 . For spin 4, its 9 components are distributed over A_1 , E , T_1 , and T_2 .

To identify to which representation an operator belongs requires to identify how it transforms under the discrete subgroup. E. g. an operator invariant under rotations will belong to the representation A_1 , while a (discrete) vector will transform as an element of the T_1 representation. It is thereby possible to construct systematically the assignment of a spin to a given operator²⁴. However, as the five representations are independent, continuum degeneracies do not persist on the lattice, except within a representation. The two and three spin states of spin 2 belonging to E and T_1 would be degenerate within a representation, but not with each other, $E_E \neq E_{T_2}$. This degeneracy will only be recovered in the continuum limit, when the rotation group reemerges.

The difference occurring between A_1 and A_2 , and T_1 and T_2 have to do with the spinor representations of $SO(d)$, which are given by Dirac matrices. The Dirac matrices 1, γ_0 , γ_5 , and $\gamma_0\gamma_5$ belong to A_1 , while γ_i , $\gamma_0\gamma_i$, $\gamma_5\gamma_i$, and $\gamma_0\gamma_5\gamma_i$ all belong to T_1 , forming together the 16-dimensional spinor basis.

These issues become quickly rather complex if considering spins greater than one. It should also be noted that the assignment of (charge)parity has to be taken care of, though

²⁴In practice, this can become quite tedious. See, e. g., arxiv:0803.4141.

as a discrete symmetry it usually just adds to the corresponding representation, giving then R^{PC} assignments as usual, but R is now one of the discrete representations, rather than the continuum spin. Note that positive and negative parity do not necessarily have to be connected in the same way to the various discrete representations.

2.9.2 Internal symmetries

Genuine internal symmetries are not affected by the presence of a lattice regulator. E. g., upgrading the ϕ^4 model to an $O(N)$ requires only the same steps as in the continuum, i. e. adding further indices to the fields, and sorting quantities by their respective charges. This is also true for local symmetries like gauge symmetries.

Especially, this implies that a lattice formulation, viewed as a ultraviolet or infrared regulator, does not break such symmetries. This is not necessarily the case, e. g. cutoff regulators break the internal symmetries, especially gauge symmetries. This is, however, not a conceptual problem, as this is similar to the breaking of the space-time symmetry by a lattice. When removing the regulator consistently, all symmetries are recovered.

Using the lattice regulator is thus nothing but a different choice of which symmetries are broken by the regulator. In fact, it is an unfortunate consequence of our current understanding of quantum field theory that the necessity of having regulators always imply that some symmetries are broken, no matter what, until the regulator is removed.

Thus, lattice is the choice of regulator if to keep internal symmetries intact when regulating is of paramount importance²⁵.

The previous statements are no longer true if the internal symmetries and space-time symmetries are intertwined, as happens for supersymmetry and chiral (gauge) theories, not to mention general relativity. In this case also the internal symmetries are broken, inducing differing amounts of problems. These problems range from inconvenient to making a lattice regularization unusable, even if they are recovered in the continuum limit. These issues will be discussed in more details in chapter 6.

There is also an issue with the so-called spontaneous breaking of symmetries. A lattice formulation is a consistent non-perturbative formulation. In such a formulation, a thing like spontaneous symmetry breaking does not occur. The reason is that the usual way of defining spontaneous symmetry breaking requires to make a choice of preference. This choice is always human-made, and therefore does not happen on its own. Thus, all order parameters always vanish.

²⁵Similar properties are often associated with dimensional regularization. However, this is not true beyond perturbation theory, as dimensional regularization requires analyticity properties of correlation functions, which are most likely not fulfilled beyond perturbation theory.

What happens is that the theory remains in a metastable state, which would show a breaking of the symmetries in case of any external perturbation, no matter how weak. In absence of such an external perturbation, this will not happen. This is actually not a feature of a lattice formulation, but is true in any genuine non-perturbative formulation. But it is something encountered somewhat regularly in the context of numerical simulations.

How to deal with this problem numerically will be described in section 4.8. There are two possibilities how to deal with the situation conceptually.

One is to use only observables, which respect the symmetry, and will therefore not vanish. There are observables of this type, which can still be used to test for the metastability. The drawback is that such operators usually involve more individual field operators than, say, order parameters. Since numerical noise quickly increases with the number of operators, often exponential or worse, this is not always possible.

The alternative is that an external source is introduced to explicitly break the symmetry. The simulations are done for several values of the external source and an extrapolation to zero is performed afterwards. However, this limit is not for all quantities analytic, as the theory itself is not analytic in the limit, and therefore may be problematic. This is even more troublesome in numerical calculations, where non-analyticities are notoriously hard to detect.

Chapter 3

Measurements

After having now a basic formulation of a lattice theory at hand, the next important step is to understand how to obtain results from it. At first this seems trivial, as, after all, this would be just evaluating the expectation values $\langle \mathcal{O} \rangle$. And this seems to be quite similar as in the continuum, as the case of perturbation theory in section 2.6 seems to have shown. While this is true, and all statements of the continuum theories remain true, though perhaps somewhat more tedious, there are now additional possibilities, which are usually not considered in the continuum. Probably one of the most important ones is that operators on a finite lattice are not distributions in space-time, but rather ordinary functions. They can therefore be used in quite different ways. Also, the formulation in Euclidean space-time allows somewhat different opportunities.

3.1 Expectation values

The first insight is how to consider expectation values. In the continuum, these expectation values are usually considered to be some function, which fulfills certain (self-consistency) equations. On a lattice, an expectation value can be rather considered as its original definition of an average.

Consider a lattice and a single scalar field. By giving the finite set of numbers $\{\phi(x_i)\}$, where x_i are all possible points on the lattice, a full space-time history of the system is completely described. Such a set of values is called a configuration \mathcal{C} . Expectation values can thus be considered as

$$\langle \mathcal{O} \rangle = \frac{1}{N_{\mathcal{C}}} \sum_{\mathcal{C}} e^{-S} \mathcal{O} = \frac{1}{N_{\mathcal{C}}} \sum_{w\mathcal{C}} \mathcal{O} \quad (3.1)$$

where in the second step the w indicates that the summation needs to be weighted by the Boltzmann factor e^{-S} , to make the nature of the summation as a sum over configurations

really explicit. The number N_C is the number of configurations included, which, in principle should be all. This will become particularly important in chapter 4.

If the field would be restricted to a finite number of values, e. g. if ϕ would be a spin variable, then this sum would indeed be finite: Since there is a finite number of lattice points and the field can only take a finite number of values on every lattice point, the sum is indeed a sum, and thus has a finite number of terms on any finite lattice. If the field can take a value in a continuous range, like the field of the usual ϕ^4 theory, the sum is not a real sum, but needs to integrate over all field values. However, as long as the lattice is finite, it is still possible to exchange the summation over configurations with such integrals (if technically possible), an important difference compared to the continuum, and which is one of the reasons why it is possible to make very powerful mathematical statements on finite lattices. When taking the infinite-volume limit, this is no longer true, and care has to be taken.

It is this view that expectation values are (weighted) sums over configurations, which lies at the heart of most investigations using lattice methods.

A consequence of this very literal implementation of expectation values of a path integral (2.6) is that the correlation functions so obtained are always the full, non-amputated ones.

3.2 Spectroscopy and energy levels

Euclidean space-time introduces a number of interesting consequences. Consider an arbitrary operator. Work in the Heisenberg picture which puts all the time-evolution in the operator. The time-evolution operator is now given by $\exp(\pm Ht)$ for a time-independent Hamilton operator. Then an expectation value can be written as

$$\begin{aligned} \langle \mathcal{O} \rangle &= \langle e^{Ht} \mathcal{O} e^{-Ht} \rangle = \sum_{nm} \langle 0 | e^{Ht} | n \rangle \langle n | \mathcal{O} | m \rangle \langle m | e^{-Ht} | 0 \rangle \\ &= \sum_{nm} e^{(E_n - E_m)t} \langle 0 | n \rangle \langle m | 0 \rangle \langle n | \mathcal{O} | m \rangle. \end{aligned}$$

Thus, various overlaps appear and an exponential damping factor.

Consider now the case of a particle created by an operator at time t_1 and destroyed at time t_2 , suppressing all other arguments. Then

$$\begin{aligned} \langle \mathcal{O}(t_2)^\dagger \mathcal{O}(t_1) \rangle &= \langle \mathcal{O}(0) e^{-H(t_2 - t_1)} \mathcal{O}(0) \rangle = \sum_n \langle \mathcal{O}(0) e^{-H(t_2 - t_1)} | n \rangle \langle n | \mathcal{O}(0) \rangle \\ &= \sum_n e^{-E_n(t_2 - t_1)} |\langle \mathcal{O} | n \rangle|^2 \end{aligned}$$

where it has been used in the first step that $H|0\rangle = 0$. The matrix elements $\langle\mathcal{O}|n\rangle$ are called the overlap of the operator \mathcal{O} with the state n . They will depend on any other quantum numbers or parameters besides the energy, as will the energies themselves.

It should be noted that it was assumed that the overlap norms are positive (semi-definite). While this is true in non-gauge theories, and in gauge theories for gauge-invariant observables, this is not generically true. Especially gauge-dependent quantities, say the photon or electron propagator, have generically negative norm contributions. This complication will be ignored for the moment, until section 5.5. Thus for now it can be assumed that all correlation functions have this form.

Consider now the special case $t_1 = 0$, then

$$\langle\mathcal{O}(t)^\dagger\mathcal{O}(0)\rangle = \sum_n e^{-E_n t} |\langle\mathcal{O}|n\rangle|^2 \stackrel{t \gg E_n^{-1}}{\approx}_{n>0} |\langle\mathcal{O}|0\rangle|^2 e^{-E_0 t}. \quad (3.2)$$

Thus, at sufficiently long time (separations) any expectation value will be dominated by the ground state, a consequence of the Euclidean space-time. At least, as long as the overlap with the ground-state is non-zero. Otherwise, the lowest level with non-zero overlap will dominate.

This result shows two important features.

Any correlation function is characterized by two sets of numbers. One are the energy levels E_n for the given quantum number channel. They will be the same for all correlation functions in a fixed quantum number channel. The other are the overlaps of the operator \mathcal{O} defining the correlation function with the energy eigenstates $|n\rangle$. These are characteristic for the correlation function in question, sometimes called the fingerprint of the operator.

The second feature comes from the fact that the energy levels in a finite volume are discrete, and therefore there is a denumerable infinite number of them. Thus, there exists a linear transformation, a unitary matrix O , between any complete base of operators and the energy eigenbasis. Conversely, knowing the expansion (3.2), it is possible to reconstruct this transformation. This feature is used in the so-called variational analysis, to be discussed in section 4.6.

It should be noted that if one of the involved operators has a vacuum expectation value, i. e. $\langle\mathcal{O}\rangle \neq 0$, there are disconnected contributions. These have to be subtracted, to obtain the connected correlation functions,

$$\langle\mathcal{O}(0)\mathcal{O}(t)\rangle_C = \langle\mathcal{O}(0)\mathcal{O}(t)\rangle - \langle\mathcal{O}(0)\rangle^2,$$

as otherwise the vacuum expectation value will contribute a constant in (3.2). While conceptually not an issue, it is in numerical calculations a problem due to the degrading signal-to-noise ratio at long times, as discussed in section 4.4.

3.3 Bound states

So far, the operators had been any arbitrary ones. However, the previous discussion becomes extremely useful in the case of states creating particles. Consider some given quantum number channel identifying a particle. In the ϕ^4 theory this could be a scalar particle, which would then be just the field ϕ . However, this does not need to be necessarily a single field operator, and any composite operator with the same quantum numbers would serve the purpose as well, e. g. $\phi(x)\phi(x) = (\phi\phi)(x)$. This will become very important in gauge theories in chapter 5, as only composite operators can be gauge-invariant.

Select now the rest-frame of the generated particle, and perform a Fourier transform on the spatial momenta

$$\mathcal{O}(t, \vec{p}) = \sum_x e^{i\vec{p}\vec{x}} \mathcal{O}(x) \stackrel{\vec{p}=\vec{0}}{=} \sum_{\vec{x}} \mathcal{O}(t, \vec{x}) = \mathcal{O}(t), \quad (3.3)$$

for zero spatial momentum, giving an operator only dependent on the time. This is sometimes called a projection on zero momentum.

However, in the rest frame, the energy is just the mass, and thus

$$\langle \mathcal{O}(t) \mathcal{O}(0) \rangle = \sum_n e^{-E_n t} |\langle \mathcal{O} | n \rangle|^2 \stackrel{t \gg E_n^{-1}}{\approx} |\langle \mathcal{O} | 0 \rangle|^2 e^{-Mt} \quad (3.4)$$

Thus, at sufficiently long times the mass of the ground state particle can be determined from the exponential decay of the correlator.

If there are additional stable bound states in the same channel, all energy levels up to the elastic threshold will be again the masses of these bound states. Thus the spectrum of stable states can be determined from the lowest exponential fall-offs.

There is, however, a complication, in interpreting these masses directly as the desired quantities. As was already visible in section 2.6.1, already perturbatively the masses on the lattice will usually not agree with the continuum ones. In particular, the mass will depend on the volume. This can be immediately seen by the fact that the size of state is affected by quantum fluctuations, and the virtual particle cloud around it has a finite extent. If this extent is of similar size as the lattice volume, the particle will be distorted.

For a stable particle, this distortion can be perturbatively estimated. Neglecting discretization artifacts, this leads for the ϕ^3 theory to

$$m = M + \frac{3}{16\pi^2 m(aN)} \left(\lambda^2 e^{-\frac{\sqrt{3}}{2} m(aN)} + \frac{m}{\pi} \int_{-\infty}^{\infty} dy e^{-\sqrt{m^2 + y^2}(aN)} F(iy) + \mathcal{O}\left(e^{-m'(aN)}\right) \right). \quad (3.5)$$

Herein is F the forward scattering amplitude, m is the infinite-volume mass of the state, and m' the next higher mass appearing in the corresponding quantum number channel. While non-trivial in itself, this correction is not larger than the leading exponential term. Thus, volume corrections to the mass of a stable state fall off exponentially with volume. While the particular form is special to the ϕ^3 theory, the statement remains generally true: Stable particle masses approach their infinite-volume limit exponentially fast.

Likewise, it can be shown that the mass of a stable bound state of two particles of mass m behaves as

$$m_B = M_B - \frac{3g^2}{16\pi m_B^2(aN)} e^{-\sqrt{m^2 - \frac{m_B^2}{4}}(aN)} + \mathcal{O}\left(e^{-m'_B(aN)}\right) \quad (3.6)$$

where g is the coupling of the interaction binding the particles together.

Thus, the energy levels can not be directly interpreted as the infinite-volume masses, but approach them exponentially fast at sufficiently fine discretizations. Physically, this can be understood, because the propagation for a massive particle is exponentially damped in Euclidean space-time. Since the finiteness of the box is accompanied by periodic boundary conditions¹, the distortion comes essentially from self-interference of a particle moving around the box. This becomes then exponentially damped.

The situation is drastically different if the particles are massless. Then their propagation is no longer damped exponentially, but only by inverse powers of the box size. However, because of the finite volume, such a state will have a mass², as the volume acts as an infrared regulator, of order $1/(aN)$. As a consequence, the mass $M \sim 1/(aN)$, and thus the leading terms in both (3.5) and (3.6), are polynomial in the lattice size, giving the leading contribution as a finite-volume artifact for massless particles.

The situation changes, if the elastic threshold is crossed, and the states become resonances. Of course, there exist channels in which no stable particles exist, and then also these considerations do not even apply for the ground state. This is the next topic.

3.4 Resonances, scattering states, and the Lüscher formalism

Besides bound states there are two other categories which play a role for spectroscopy: Scattering states and resonances.

¹Other boundary conditions yield different particularities of the effect, but qualitatively yield also finite-volume corrections to masses.

²For particles with spin one or larger, this becomes more involved.

Scattering states are states which involve multiple individual particles, which may or may not have relative momentum, and form a given set of total quantum numbers. These particles are considered to be essentially non-interacting, i. e. sufficiently far separated to have at most exponentially small interactions³. In case of the ϕ^4 theory for the 0^+ quantum number channel a single-particle state is $\phi(t, \vec{0})$. The two lowest scattering states in the same channel are $\phi(t, \vec{0})\phi(t, \vec{0})$ and $\phi(t, \vec{p})\phi(t, -\vec{p})$, where \vec{p} is the smallest lattice momentum, i. e. a momentum with a single non-zero component of size π/N .

Considering the non-interacting case, the masses of these state are, up to lattice corrections, m , $m_s^0 = 2m$ and⁴ $m_s^1 = 2\sqrt{m^2 + \vec{p}^2}$. The first scattering state is then the lightest state into which a heavy scalar particle of mass $3m > M > 2m$ could decay, it gives the elastic threshold. Likewise, the scattering state $\phi(t, \vec{0})^3$ defines the inelastic threshold for a particle of mass $M > 3m$, as such a particle can now decay into either two or three particles.

As a consequence, in every quantum number channel there are always all possible scattering states. Scattering states of massive particles without internal relative momenta receive only the exponential corrections due to finite volume. However, states with relative momenta receive polynomial corrections, as the momenta depend polynomially on the volume. This can actually be helpful, as it pushes scattering states with non-zero relative momenta high up in the spectrum on sufficiently small volumes, and there is hence a sweet spot in volume where the exponential corrections are yet small but the polynomial ones are still large.

It should be noted that relative momenta can alter the quantum number of states, as momentum carries an angular momentum of 1^- .

The other category of possible states are resonances, i. e. unstable bound states. Resonances are in itself a quite intricate topic in quantum field theory. They show up as additional states in the energy levels above the elastic threshold. Thus, their presence can be detected by level counting. Surplus states compared to the expected scattering states are derived from resonances. Unfortunately, their properties are much harder to determine.

To see this note that a finite lattice is essentially quantum mechanics. Thus, there is avoided level crossing, i. e. if the system is deformed by a change of the parameters, in particular the volume, which affects the energy levels, they will continuously change, but

³Interactions by massless particles pose a problem here, as they make strictly speaking this idealization impossible. However, aside from QED all currently relevant theories on the lattice have a mass gap, and therefore the following applies.

⁴To include discretization artifacts the formula $2 \cosh(m_s^1 a/2) - 2 \cosh(am) - (2 \sin(ap/2))^2 = 0$ can be used. However, this correction is usually quite small for small momenta.

never cross. As resonances are above the elastic threshold, the number of energy levels below them will continuously increase up to infinity when sending the volume to infinity. Thus, every time a scattering state would cross the resonance level actually the level do not cross, but swap their identity instead. Moreover, the energies are real. This is fine for stable states, which are associated with real poles in Minkowski space-time, but not for resonances, which belong to complex poles on the second Riemann sheet. Thus, the energy levels give only a summary information, which may or may not be similar to the real part, and thus mass, of the pole. Thus, at first sight it seems hopeless to learn about resonances from lattice field theory. Of course, the information is, due to the reconstruction theorem, contained in the theory. But this leaves the question how to obtain it. This is even more serious in practical calculations, where only a finite number of parameter sets and lattice points are available, which cannot determine an analytic function as a correlation function uniquely.

Fortunately, there is a way around it.

The basic physics idea is that two particles in a finite box will always interact. This interaction will distort the energy levels of this scattering state. It is then possible to derive from these deviations information about this interaction. Since this is the same type of interaction as responsible for the binding of the resonances, this implies that it should be possible to infer properties of the binding mechanism, and thus resonance properties, from these distortions. This is indeed possible using the so-called Lüscher method. Like for the determination of the volume-dependence of the masses it uses essentially a perturbative analysis to determine how the volume, assuming again negligible lattice spacing effects, distorts the energy levels.

Consider for now the case of a resonance at rest decaying into two identical, spinless particles of mass m , and that this is the only open decay channel. The method has been extended to include different particles, also with spin and outside the rest frame of the resonance, and to include multiple open two-body channels, and the development of an extension to three-body channels is essentially finished at the time of writing. While technically far more complicated than the simple case they are conceptually quite similar. Therefore, here only this simplest case will be discussed. If need be, the other cases can yet only be found in contemporary literature.

As deformations of the energy levels occur due to interaction, and thus scattering, of the constituents, the central quantity of interest are the phase shifts. The elastic cross-section, i. e. precisely the process which a scattering state undergoes, reads

$$\sigma_e = \frac{16\pi^3}{s} \sum_l (2l + 1) \left| \frac{e^{2i\delta_l(s)} - 1}{2i} \right|^2,$$

where s is the center-of-mass energy, and the sum is over all partial waves of angular momentum l . The quantities δ_l are then the phase shifts, which characterize the elastic process completely. The phase shifts in turn are connected to the width function Γ as

$$\cot \delta_l(s) = \frac{M^2 - s}{\Gamma(s)\sqrt{s}}$$

where M is the mass of the resonance of angular momentum l . If the resonance has a pole close to the real axis (in comparison to its mass), the amplitude has the form.

$$A(s) \sim \frac{-\sqrt{s}\Gamma(s)}{s - M^2 + i\sqrt{s}\Gamma(s)} = e^{i\delta_l} \sin(\delta_l) \quad (3.7)$$

The cross-section is then a Breit-Wigner one

$$\sigma \sim \frac{\Gamma^2}{(\sqrt{s} - M)^2 + \frac{\Gamma^2}{4}}$$

where Γ is evaluate at M^2 . Wide resonances show a quite different behavior.

Now, there will be energy levels E_n between the elastic and inelastic threshold on the lattice. These energy levels are determined as

$$E_n = 2\sqrt{m^2 + \vec{p}_n^2}.$$

Without interactions, \vec{p}^2 would be just the normal momenta, i. e. given by integer multiples of $2\pi/N$. Now, they will be given by⁵

$$\vec{p}_n^2 = \left(q_n \frac{2\pi}{N} \right)^2$$

with some, usually non-integer, number q_n . This number can be determined once the energy level E_n and the infinite-volume masses m of the decay products are known.

It has then been shown that these numbers are related to the scattering phase shifts, where the exact form depends on the relative momentum and the spins of the involved particles. In the simplest case of spinless particles with vanishing center-of-mass momentum this relation reads

$$\tan \delta_l(q_n) = \frac{\pi^{\frac{3}{2}} q_n}{Z_{00}^0(1, q^2)}.$$

⁵Lattice artifacts could be included in this relation e. g. by using the improved dispersion relation of footnote 4.

The function $Z_{00}^{\vec{0}}$ is a purely geometric function, which is obtained from analytically continuing

$$Z_{lm}^{\vec{d}}(r, q^2) = \sum_{\vec{x} \in P_{\vec{d}}} \frac{|\vec{x}|^l Y_{lm}(\vec{x})}{(\vec{x}^2 - q^2)^r}$$

$$P_{\vec{d}} = \left\{ \left(\vec{m} + \frac{\vec{d}}{2} \right) \mid \vec{m} \in \mathbb{Z}^3 \right\}.$$

The functions Y_{lm} are the usual spherical harmonics. This sum is not convergent for $r = 1$, as it is essentially the ζ -function, and therefore needs to be analytically continued. A possible continuation yields at $r = 1$

$$Z_{lm}^{\vec{d}}(1, q^2) = \int_0^1 dt \left(\frac{\pi}{t} \right)^{\frac{3}{2}} e^{tq^2} \sum_{\vec{u} \in \mathbb{Z}^3, \vec{u} \neq 0} (-1)^{\vec{u}\vec{d}_i l} \left(-\frac{\pi|\vec{u}|}{t} \right)^l Y_{lm} \left(-\frac{\pi|\vec{u}|}{t} \right) e^{-\frac{\pi^2|\vec{u}|}{t}}$$

$$+ \frac{\delta_{l0}\delta_{m0}}{\sqrt{4\pi}} \int_0^1 dt \left(\frac{\pi}{t} \right)^{\frac{3}{2}} (e^{tq^2} - 1) - \pi\delta_{l0}\delta_{m0} + \sum_{\vec{x} \in P_{\vec{d}}} \frac{|\vec{x}|^l Y_{lm}(\vec{x})}{\vec{x}^2 - q^2} e^{-(\vec{x}^2 - q^2)},$$

which converges exponentially fast.

Thus, it is possible to reconstruct the phase shifts from the energy levels. However, in any practical calculations usually only very few energy levels are located between the elastic and inelastic threshold, sometimes not more than one. Thus, with a single lattice it is usually not possible to obtain enough information to reconstruct the phase shift. This can be improved by either using several different volumes or by adding more energy levels. The latter is obtained by also including moving scattering states, i. e. with non-zero center-of-mass momentum, or, equivalently, a resonance not decaying at rest. This latter option complicates things further. It is then necessary to transform the result back into the rest frame. This is straightforward, but cumbersome in detail.

Combining the the result for the phase shift with the amplitude relation (3.7) yields

$$\sqrt{s}\Gamma(s) \cot \delta_l(s) = M^2 - s$$

and thus the width can be obtained from this relation. As the width itself includes phase space effects, it is often convenient to exchange it for the coupling $g(s)$ of the decay channel to the resonance by

$$\Gamma(s) = \frac{p^3 g(s)^2}{2 \cdot 6\pi s}$$

which can also be determined directly from the phase shifts

$$\frac{p^3}{\sqrt{s}} \cot(\delta_l(s)) = \frac{6\pi}{g(s)^2(M^2 - s)} \approx a + bp^2.$$

where the linear expansion in p^2 holds only close to zero momentum of the constituents. Herein p is the would-be three-momentum the scattering states would have at the given s , and will coincide with the $\vec{p}_n/2$ for the corresponding values of s which agree with measured lattice energies. The parameter a is known as the scattering length. As this quantity is uniquely connected to the coupling, it could, in principle, be also extracted directly from the finite-volume modification of the bound state (3.6). However, in practice the method presented here is superior, and forms the standard approach.

3.5 General correlation functions

While masses and widths are usually the primary quantities of interest, they are not always sufficient. An alternative are fully resolved correlation functions in momentum space⁶.

Calculations of such correlation functions proceed in a rather straight-forward way. Given any field, its full momentum-space version is obtained by a discrete Fourier transformation

$$\phi(p) = \frac{1}{\sqrt{\prod_{\mu} N_{\mu}}} \sum_x e^{i2\pi \sum_{\mu} \frac{p_{\mu} x_{\mu}}{N_{\mu}}} \phi(x).$$

The p_{μ} are here the integer-valued lattice momenta. Note that the Fourier-transformed of even a purely real field is usually complex. Note that the complex conjugate yields

$$\phi(p)^{\dagger} = \frac{1}{\sqrt{\prod_{\mu} N_{\mu}}} \sum_x e^{-i2\pi \sum_{\mu} \frac{p_{\mu} x_{\mu}}{N_{\mu}}} \phi(x)^{\dagger}.$$

Thus, for a purely real field, the complex conjugate in momentum space only reverses the momentum. This implies that only momenta with $0 \leq p_{\mu} \leq N_{\mu}/2$ provide independent information, due to the periodic boundary conditions.

Arbitrary correlation functions can then be obtained by calculating products of the fields, and then averaging them over configurations. E. g. a propagator is given by

$$D = \langle \phi(p)^{\dagger} \phi(p) \rangle = \frac{1}{N_c} \sum_i \phi_i(p)^{\dagger} \phi_i(p), \quad (3.8)$$

where i counts the configurations and N_c is the number of configurations, where eventually $N_c \rightarrow \infty$ should be taken, as noted already in section 3.1. For a vertex, this yields

$$\Gamma = \langle \phi(p) \phi(q) \phi(k) \phi(-p - q - k) \rangle = \frac{1}{N_c} \sum_i \phi_i(p) \phi_i(q) \phi_i(k) \phi_i(-p - q - k). \quad (3.9)$$

⁶On a finite lattice configuration-space quantities are also well defined. However, in the continuum limit, they generically become tempered distributions, and are therefore more involved. Thus, momentum-space quantities, which are essentially ordinary functions in Euclidean space-time, are usually better suited.

There are, however, a number of caveats.

The first, and probably simpler one is that these are lattice quantities, and therefore cannot be directly compared with the continuum version. E. g., to obtain a more continuum-like version of the propagator, it would be necessary to calculate

$$D(P) = a^2 \kappa \langle \phi^\dagger \phi \rangle \left(\frac{2}{a} \sqrt{\sum_{\mu} \sin^2 \left(\frac{\pi p_{\mu}}{N_{\mu}} \right)} \right)$$

where the improved lattice-momentum (2.16) is used to determine the momenta. The factor κ appears because of the rescaling of the the fields (2.20). While the lattice ensures that the so obtained propagator is regulated, and therefore finite, it is not yet renormalized. This manifests in a dependency on a , usually making the propagator either vanish or diverge for all momenta as $a \rightarrow 0$. It thus needs to be renormalized to become meaningful. But this can be done in exactly in the same way as in continuum theories, and will therefore not be detailed here.

Furthermore, if the fields carry indices, e. g. when there would be two scalar fields ϕ^1 and ϕ^2 , the propagator becomes matrix-valued

$$D^{ij} = \langle \phi^i(p)^\dagger \phi^j(p) \rangle$$

Just as in the continuum the propagator can then be decomposed into invariant tensors to isolate dressing functions (or form factors). If there is a source breaking the symmetry explicitly as

$$D^{ij} = (\delta^{ij} - n^i n^j) D_s + n^i n^j D_t$$

with the two independent scalar dressing functions D_s and D_t and n^i is a normalized vector in the direction of the source. Of course, any choice of basis is equally well justified. The dressing functions can then be isolated by projections

$$\begin{aligned} D_s &= \frac{1}{2} (\delta^{ij} - n^i n^j) D^{ij} \\ D_t &= n^i n^j D^{ij}, \end{aligned}$$

as is also done in the continuum. If the symmetry is intact, there is only one dressing function,

$$D^{ij} = \delta_{ij} D.$$

In general, the decomposition has to be done in terms of the invariant tensors of suitable rank for the given correlation function and symmetries of the involved fields, according to the Wigner-Eckart theorem.

The situation only becomes somewhat involved if the indices are not internal, but space-time, i. e. Lorentz indices. Because of the violation of rotational symmetry additional tensor structures can arise in comparison to the continuum. These vanish when the continuum limit is approached, but will at finite lattice spacing contribute. This is of special importance when lattice perturbation theory is used to estimate discretization artifacts, and adds further complications. This holds also true for vertices.

Another problem appears when considering vertices like (3.9). The fact that the vertex function are directly created from the fields implies that this is the full correlation functions. To obtain what is usually considered to be the vertex function, the connected part of this function must be taken and amputated. This also works like in the continuum, in principle. In the practice of numerical simulations, the necessary subtractions usually incur an enormous amount of statistical noise, making an extraction difficult. This also implies the necessity to know all relevant lower n -point functions to perform this procedure. Of course, as for the propagator, lattice corrections, normalization, and the projection to dressing functions have to be performed in essentially the same way as for the propagator.

The second issue has far more implications. An expression like $\phi_i(p)^\dagger \phi_i(p)$ in (3.8) hides the important fact that on a single configuration no symmetries exist. Especially, any single field configuration has no rotational or translational symmetry. This implies that on a single configuration there is also no momentum conservation, and the propagator should be actually rather $\phi_i(q)^\dagger \phi_i(p)$, and the dependency on p and q will be different. Only summation over many configurations will restore momentum conservation and a dependency purely on $p^2 = q^2$ rather than the individual components. This requires the number of configurations N_c in (3.8) and (3.9) to be large enough that these effects become small. Performing an evaluation only of $\phi_i(p)^\dagger \phi_i(p)$ is, in a sense, already an improvement for a finite number of lattice configurations.

To assess the extent of such problems it is useful to evaluate correlation functions along different momentum axes. In the continuum, they agree. On the lattice, they will in general not. Comparing results with permuted momentum components will give an estimate of how badly momentum conservation and reflection symmetries are deteriorated by a finite number N_c . Comparing results with different momentum components but equal length of the momentum gives an estimate of violation of rotational symmetry.

Finally, even for theories without interactions the propagator has not the free form (2.15) on a single configuration. In fact, even without interactions the fields $\phi(x)$, and thus $\phi(p)$, fluctuate strongly from site to site⁷. Only on the average this strong fluctuations

⁷This is what is expected from the continuum: The largest contribution to the path integral stem from non-differentiable field configurations, corresponding to lattice fields fluctuating arbitrarily between

cancel each other, giving the simple free propagator.

lattice sites. This is a very literal realization of the concept of quantum fluctuations.

Chapter 4

Monte-Carlo simulations

While the approximation of quantum field theory using quantum mechanics provides the possibility for stronger mathematical statements in some cases, the real advantage is that it makes such theories accessible to numerical simulations. In particular, such simulations are possible for many Lagrangians no matter the couplings, and therefore can cover without further approximation also non-perturbative physics. This allows to directly probe bound states and phase transitions, which are genuine non-perturbative phenomena.

The drawback is that there are various error sources, which can lead to both qualitative and quantitative problems. These will be addressed in sections 4.4 and 4.5.

4.1 Importance sampling and the Markov chain

Formally, the lattice path integral is just a large number of integrals which, in principle, could also be numerically evaluated, as can be seen from its measure (2.4). However, the number of integrals scales as N^d , and thus very quickly not all integrals can be performed individually within any reasonable amount of computing time. Thus, a different approach is necessary.

The basic idea of the solution is already given in section 3.1: Determine configurations, and then average expressions over them, appropriately weighted by the Boltzmann factor.

The practical problems are twofold. One is how to obtain configurations. The other how to obtain enough configurations.

The first answer seems to be simple, as any values for the field variables is actually a valid field configuration. This, however, leads to a serious problem with the second part of the question. Some arbitrary field configuration will usually have a very small Boltzmann factor. Only rarely there will be a configuration with a large Boltzmann factor, and the ratio will worsen with volume. The reason is that most configurations are more or less

just quantum noise, which is absorbed during renormalization. The relevant information is hidden in this noise.

The solution to this problem comes from the analogy of section 2.2 of a lattice theory with a stochastic system. Introducing a fictitious fifth coordinate, the so-called Monte-Carlo time, a choice of a relevant configuration can be considered as non-equilibrium equilibration. It is this idea which will now be used.

For the following, it is convenient to write the action as

$$S = \beta V \epsilon = \beta E \quad (4.1)$$

where $V = \Pi_\mu N_\mu$ is again the volume and β some factor. In case of ϕ^4 this would be κ , but it can have any definition. However, it is not by chance that it is the same symbol as the inverse temperature in statistical mechanics, and E is then the energy, and ϵ the energy density. This makes the analogy to a four(five, when including the time when out of equilibrium)-dimensional statistical mechanics system in equilibrium evident. It is also useful to define the density of states d as

$$d(\epsilon) = e^{Vs(\epsilon)} = \int \mathcal{D}\phi \delta\left(\epsilon - \frac{S}{V\beta}\right),$$

where s is thus the same as the entropy density, and ϕ are the fields of the theory, which are not necessary to specify in detail for most of the following.

This allows also to rewrite the partition function as

$$Z(\beta) = \int \mathcal{D}\phi d(\epsilon) e^{-V\beta\epsilon} = \int d\epsilon e^{V(s-\beta\epsilon)} = \int d\epsilon e^{-Vf},$$

introducing the free energy density f . In analogy to statistical mechanics

$$\rho(\beta, \epsilon) = \frac{d}{Z} e^{-V\beta\epsilon} = \frac{e^{-Vf}}{Z}$$

is then the probability density for a given energy density.

Now, in the sum over configurations only those will substantially contribute which have a not too small Boltzmann factor $\exp(-S)$. In fact, if the system is viewed as a statistical ensemble, then any equilibrium configuration has a weight, or density, W of

$$W[\phi] \sim e^{-S[\phi]}. \quad (4.2)$$

The constant of proportionality will be (implicitly) fixed such that

$$\sum_{\phi} W[\phi] = 1,$$

i. e. the total weight of all configurations is unity. If it is now possible to create configurations with a probability according to this weight, then for any subsample of N configurations any average

$$\langle A \rangle = \frac{1}{N} \sum_{\mathcal{C}} A[\mathcal{C}] \quad (4.3)$$

will approach the same average as a sum over all configurations with the weight functions as in¹ (3.1). However, outside of equilibrium the configurations will have a differing weight W_n . In this case, (4.3) will not hold, as this corresponds to taking the path integral with a different weight function.

Note that though written like an expectation value for any finite number of configurations N this number will be different from the actual expectation value. This is a statistical error, which will be discussed in more detail in section 4.4. This is sometimes emphasized by using different symbols for the case of finite and infinite (all) N .

This leaves the question how to actually obtain such a selection of weighted configurations.

In practice this is solved by a method which deforms a given configuration ϕ_n into a new configuration ϕ_{n+1} by a sequence of changes, which are made such as to make a configuration relevant. This is called an update. This seems to require that a first good configuration exists. However, a suitable chosen update process will move any configuration towards a relevant one by updates. Thus, in practice a simulation starts by some, often random, field configuration, which is then made relevant. This is called thermalization. This will be discussed more in section 4.4.2. For the moment, just consider the situation of getting from one configuration to a new one.

It is, in principle, irrelevant which subset of configurations is used in (4.3), and in fact it is best to have a random sequence, again for reasons to be discussed in section 4.4.2. Thus, the transition between two configurations can be considered a stochastic process, and a transition from a configuration ϕ to a new configuration ϕ' will occur with a probability $P(\phi \rightarrow \phi')$. This probability is normalized

$$\sum_{\phi'} P(\phi \rightarrow \phi') = 1$$

i. e. given any configuration the total probability for a transition to any other configuration is one: Configurations are not lost in transition. Also, the probability needs to be positive for any configuration ϕ to give a reasonable probability². This is also called strong

¹This implies that in the limit of infinite configurations, configurations would be included multiple times in this sum, with the multiplicity given by $\exp S$. However, for any reasonable amount of computing time, no configuration will be contributing even twice.

²Which is precisely the origin of the sign problem to be discussed in section 7.2.

ergodicity, as it implies that any configuration can be reached from every configuration with a finite probability. Note that this does not require that this is a direct transition. In fact, in practice a transition is often build up from many steps, which individually do not satisfy strong ergodicity, but combining a certain number of them in a single update is strongly ergodic. E. g., most updates consist out of steps which change the configurations only at a single lattice site. Such a step is not strongly ergodic, as it cannot connect two configurations which differ at different lattice sites. However, combining steps such that every lattice site is updated can constitute a strongly ergodic update.

The transition probability is related to (4.2) as

$$W_{n'}[\phi'] = \sum_{\phi} P(\phi \rightarrow \phi') W_n[\phi]. \quad (4.4)$$

Thus, the weight of a configuration is modified by the transition to a new weight.

As noted above, it is usually not possible to guess an initial configuration such that already 4.2 holds, but it will be a different configuration. Updates will drive the configuration towards W if

$$\lim_{k \rightarrow \infty} \sum_{\phi^{i^k}} P(\phi^{i^k} \rightarrow \phi) \sum_{\phi^{i^{k-1}}} P(\phi^{i^{k-1}} \rightarrow \phi^{i^k}) \dots \sum_{\phi^{i^1}} P(\phi^{i^1} \rightarrow \phi^{i^2}) W_x[\phi^{i^1}] = W[\phi],$$

i. e. in the limit of an infinite number of transitions any weight will approach the Boltzmann weight. This is like the usual thermalization of a system out of equilibrium towards equilibrium, and the same name is used. This implies that

$$W[\phi'] = \sum_{\phi} P(\phi \rightarrow \phi') W[\phi]$$

and thus that the Boltzmann weight (4.2) is invariant under updates. Also this is expected from thermodynamics: Once a system has reached equilibrium, it will stay there forever, no matter what the individual parts of the system do. It should be noted that in case of multiple equilibrium situations, like at first order phase transitions, this requires some care, see section 4.7. If the system would have more than one equilibrium distribution, which are not related to any phase transition, it may become more complicated to give sense to the numerical process. Since this is not relevant in the systems of interest in particle physics, to the best of our current knowledge, it will be assumed not to be the case. If there is a single equilibrium situation, then actually the starting point is irrelevant, if all other conditions listed above are met.

A sufficient, but not necessary, condition for an update to satisfy all the conditions listed above is detailed balance

$$P(\phi \rightarrow \phi') W[\phi] = P(\phi' \rightarrow \phi) W[\phi'] \quad (4.5)$$

i. e. the probability for going back and fore between two configurations is given by the ratios of the Boltzmann weight. While weaker conditions can be formulated, usual lattice algorithms all fulfill this stronger condition, and this will therefore suffice here. A proof of these statements can be obtained by considering (4.4) to be a matrix-vector equation with P the matrix and W the vector. Then the fact that there is one largest eigenvalue of P which is 1, with eigenvector W , and all other eigenvalues smaller than 1 and P a positive matrix guarantees with some linear algebra theorems the rest.

Any update process which has the features above, and therefore automatically any algorithm satisfying detailed balance, is called a Markov process. The sequence of configurations created is called a Markov chain.

4.2 Metropolis algorithm

Now the requirements for an update to create configurations has been formulated, but still an explicit implementation is lacking. In fact, there exist a very large number of possible algorithms. Here, an example will be discussed, the so-called Metropolis algorithm. It is very flexible and can be used for almost any theory for which a simulation is reasonable at all. However, it is not a particular efficient one, and theory-specific ones can easily outclass this algorithm by orders of magnitude in performance. It is therefore strongly advised to thoroughly search the literature before choosing an algorithm.

To flesh out the algorithm requires essentially only one thing: To provide an update which changes a configuration, and to show that it fulfills detailed balance.

The algorithm operates as following: First, select a new configuration. This can be done in any way thought possible. E. g. for the ϕ^4 theory every field value could be added some Gaussian distributed random number with some fixed width being the same for all lattice sites. This configuration is always accepted if its action is smaller, and therefore the Boltzmann weight (4.2) larger, than for the previous configuration. If the action is increased (W decreased) accept it with the probability $\exp(S_{\text{old}} - S_{\text{new}}) < 1$. Up to normalization, this corresponds to the transition probability

$$P(\phi \rightarrow \phi') \sim (\theta(S_{\text{old}} - S_{\text{new}}) + \theta(S_{\text{new}} - S_{\text{old}})e^{S_{\text{old}} - S_{\text{new}}}). \quad (4.6)$$

The actual test, whether a proposed changed should be accepted if its action is larger can be obtained by drawing a random number between zero and one and compare it to this probability.

In practice, such a global update yields almost never an accepted change, making this highly inefficient. This can be changed by instead performing local updates, where a

change is only proposed at a single lattice site³. Such an update is not ergodic, but by performing such updates at every lattice site in turn makes the algorithm again ergodic⁴. Then the individual acceptance probabilities become large enough to efficiently find new configurations⁵.

A full update of all lattice sites is also called a sweep. Note that the movement in configuration space, as noted above, is akin to a stochastic process. Therefore, one update is often considered to be one unit of Monte-Carlo time, and the movement in sweeps is called a propagation in Monte-Carlo time. It is important to distinguish this from the time coordinate x_0 .

To modify the acceptance probability can also be achieved by modifying how a new field configuration is proposed. The more the fields change, the less likely the new configuration will be accepted. However, if the changes are too small then only small modifications of a given configuration will be explored. This can lead to serious problems, as will be discussed in section 4.5.2. Still, acceptance rates of changes at a single site between 10% and 90% are usually a good choice.

The Metropolis algorithm satisfies detailed balance, as its involves directly the weight (4.2). Inserting (4.6) into (4.5) shows this explicitly.

A generalization of the Metropolis algorithm is obtained by decomposing it into two probabilities, $P = P_A P_C$. P_C is now an arbitrary probability distribution to obtain from any given field configuration a new one. Defining

$$P_A(\text{old} \rightarrow \text{new}) \sim \left(\theta(W_{\text{old}} - W_{\text{new}}) + \theta(W_{\text{new}} - W_{\text{old}}) \frac{P_C(\text{new} \rightarrow \text{old})W_{\text{new}}}{P_C(\text{old} \rightarrow \text{new})W_{\text{old}}} \right) \quad (4.7)$$

yields a transition probability which satisfies detailed balance, which can be proven again by direct insertion. If the transition probability is independent of the direction, P_C drops out, and the previous case is recovered. E. g. if a local modification of the field is done by adding some random number to the field value, the change is independent of the direction. If the change would have different sizes depending on the previous field values, this is not the case, and (4.7) has to be used instead of (4.6).

³If the fields have more than one component, even for changes of a single component.

⁴Since particle physics theories have only next-neighbor connections in the action, there are the possibility to update either all sites in order, e. g. using the counting (2.12) (lexicographical) or first all even sites and then all odd sites of this quantity, which reduces the respective interference (checkerboard update). However, in practice this makes usually little difference. Also a purely random choice of selected lattice points is possible, and roughly equally effective.

⁵A more quantitative definition of 'efficiently enough' will be discussed in section 4.4.2.

4.3 Improving algorithms

While the Metropolis algorithm is comparatively simple to implement, it is far from the most efficient one. This will become even more clear in section 4.4.2. There are many faster algorithms. When using knowledge about the theory in question it is always possible to be even faster. Examples of this will be discussed in sections 5.4 and 6.9. Here, generic improvements will be discussed.

4.3.1 Acceleration

A straightforward improvement is obtained by using a multi-hit Metropolis algorithm. This algorithm is only a slight change to the ordinary Metropolis algorithm. It only alters the local updates, but the composition of a sweep from local updates is not changed.

The modification is that instead of making one proposed change, a fixed number N of proposals are made per lattice site to alter the field value. Only after this the next lattice site is updated. In the limit of large N this becomes a rather efficient algorithm. However, the repeated attempts also cost time, and therefore there exists a sweet spot in N , where this algorithm is particularly efficient. The reason is that less and less the result depends on previous information of the system, and therefore more differences in configurations are possible.

The number of hits N can in addition be optimized at run-time, making this an adaptive algorithm. If a particularly measure of efficiency is available, this can be checked after every sweep to adjust N . Similarly, this can also be done for the parameters of the proposed change for the fields, e. g. the maximum size of changes. It is important that this is only done after a full sweep, as otherwise the update becomes incoherent, treating different parts of a lattice differently. This is important as that the sweet spot of such numerical parameters quite strongly depend on the parameters of a theory, and good values cannot be predicted. In this way, it is possible to adjust them at run time. Note, however, that an optimal choice could also depend on the configuration, making optimization at best hard and only working on average.

4.3.2 Overrelaxation

The aim of the sampling is to cover as many relevant field configurations as possible. The previously described algorithms are canonical, as they allow to change the action. Since the action is equivalent to the energy, according to (4.1), the name arises.

However, in many cases such changes can be inefficient in sampling all of configuration

space, as they prefer areas of lower action. While detailed balance guarantees that in the limit of an infinite number of configurations everything is sampled correctly, this may not be so with a finite number of samples. Especially, if the space of all possible configurations has multiple areas of low action, which is actually quite common, the algorithm will tend to stay close to a minimum for very many updates, before moving on to another area of low action.

To avoid this, it is possible to use microcanonical updates, i. e. any kind of change to the field configuration such that the action remains constant. As this requires knowledge about the action, such so-called overrelaxation updates cannot be formulated in general. Since such an update moves along curves of constant action without preference for minima it will tend to get less stuck at a given (local) minimum.

While in the limit of an infinite number of configurations both microcanonical and canonical updates yield the same result for theory parameters in which the average of the action of the canonical updates equals the fixed value of the action in microcanonical updates, this is not very relevant in practice. In fact, often the value of the action is the observable, and unknown in advance. Thus, in practice usually a mixture of canonical and microcanonical updates is performed. Putting together a sequence of such updates then manifest a (macro)sweep, and it is sufficient that such a sweep satisfies detailed balance, which is guaranteed, if individually both do, and sometimes even if only one does.

Due to applications of this idea to the movement of spins towards equilibrium, the name of such microcanonical updates is also overrelaxation.

4.4 Statistical errors

Since a Monte Carlo chain is stochastic process, it is equivalent to an experiment with noise. Therefore, the average (4.3) will usually differ from the one in the limit of an infinite number of configurations. It is therefore important to determine the error associated with this process.

4.4.1 Determination and signal-to-noise ratio

As the process is of the same type as any statistical process, the statistical error for a quantity \mathcal{O} with average value $\langle \mathcal{O} \rangle$ is given after N measurements by

$$\sigma_{\mathcal{O}} = \sqrt{\frac{1}{N(N-1)} \sum_i (\mathcal{O}_i - \langle \mathcal{O} \rangle)^2} \quad (4.8)$$

where \mathcal{O}_i is the measurement of the quantity on the i th configuration. This gives the usual one standard deviation, and thus the infinite-configuration average will be within $\langle \mathcal{O} \rangle \pm \sigma_{\mathcal{O}}$ with a probability of roughly 67%, and correspondingly for higher multiples of $\sigma_{\mathcal{O}}$.

While this is a suitable error estimate⁶ for Gaussian distributed quantities, it is not generally a good estimate. Especially, there are often quantities, just think of a positive operator \mathcal{O}^2 , which are not Gaussian distributed.

To deal with this, other methods are available. The two most prominent ones are bootstrap and jackknife. They both operate in the same conceptual way: By obtaining a number of averages from pseudo-experiments. They then give as an error margin values $\langle \mathcal{O} \rangle_{\pm}$ such that some fraction, say again 67%, of all so generated average values fall within $\langle \mathcal{O} \rangle_{-} \leq \langle \mathcal{O} \rangle \leq \langle \mathcal{O} \rangle_{+}$, such that the number of averages on both sides is equal, i. e. 33.5%. This allows for asymmetric errors.

The average values for jackknife are obtained by eliminating from the set $\{\mathcal{O}_i\}$ one or more elements before calculation. If one is removed, this is called single-elimination jackknife. Then the average value is determined. This is repeated many times, up to and including N times if there are N configurations, to produce the distribution of averages. Bootstrap operates oppositely. It draws from the set $\{\mathcal{O}_i\}$ a fixed, but larger, number of values, and calculates then the average value. This is repeated multiple times, usually also for the number of configurations. It is allowed to also draw numbers twice. In fact, often the number drawn is larger than the available number of configurations.

It can be shown that both methods yield for Gaussian distributions of values in the limit of an infinite number of configurations the same value as the statistical error (4.8). In fact, if the values are Gaussian distributed, they approach the value quickly. The major advantage is to deal with non-Gaussian, in particularly asymmetric, errors. Especially jackknife has become a quasi-standard, but bootstrap is conceptually not weaker. To test how Gaussian a result is, and whether the number of configurations can be considered large, the statistical error (4.8) and the one from jackknife or bootstrap can be compared.

The mathematical theory of stochastic error estimation is very extensive, but in most cases the above is sufficient.

An important issue in statistical errors is the signal-to-noise ratio, i. e. the ratio $\sigma_{\mathcal{O}}/\langle \mathcal{O} \rangle$, the so-called signal-to-noise ratio. If this ratio is one or even larger, the operator is zero within errors. In fact, to obtain a reliable statement requires that this ratio

⁶Any so calculated results can only be an estimate, as there could be a 'unlucky' sequence \mathcal{O}_i , making the error extremely overestimated or underestimated. This problem is discussed in many details in the corresponding literature. For practical purposes, these subtleties are usually glossed over.

should be much smaller than one. A usually unambiguously accepted level for a signal is three for evidence and five for a genuine signal. In practical cases it is often not possible to reach such a large ratio for many quantities of interest. This is amplified by (3.2) in spectroscopy: It is needed to follow a signal until only the ground state contributes. At the same time, the signal drops exponentially. Thus, the signal-to-noise-ratio deteriorates exponentially. There are many possibilities to counteract this by making the prefactors of the noise small, or to some extent even the prefactor in the exponent, as described in section 4.6, but eventually always the noise wins. This becomes more troublesome the heavier the mass. In fact, in many cases it is possible to estimate the statistical noise to rise like $\exp(a_{\mathcal{O}}m_{\mathcal{O}}/m_l)$, where $m_{\mathcal{O}}$ is the relevant mass scale of the observable, a is some constant, typically specific to the observable, and m_l is usually the lightest mass. Thus, the heavier the particle, the larger the signal-to-noise problem.

The signal-to-noise ratio is a function, and can depend on the parameters of the measurement. Therefore, if there are disconnected vacuum expectation values, as discussed in section 3.2, the error of the vacuum expectation value will be independent of time. Thus, the exponentially decreasing signal has to be extracted from this background, and thus the signal-to-noise ratio degrades exponentially fast. In channels without vacuum expectation values, this is not an issue, if the number of configurations becomes large enough.

4.4.2 Autocorrelations, thermalization, and critical slowing down

An important constraint on the previous discussion of statistical errors is the implicit assumption that every measurement is independent. For this to be true, every configuration needs to be independent. This is, however, in general not true if the configurations are created by a finite number of sweeps in a Markov chain. As already (4.6) shows the changes depend on the previous configurations. Though there are possibilities to reduce this dependence, overrelaxation being one of them and another one will be discussed in section 5.4, in practical cases for particle physics theories it can never be eliminated.

As a consequence, the configurations are said to be correlated to some degree. If they are correlated, observables measured are not independent, and the fluctuations between configurations are smaller than without correlations. Thus, the actual statistical error are larger than those calculated using the methods in section 4.8. The only question is how large.

To obtain an estimate for the correlation of configurations the so-called autocorrelation time can be determined, where time refers here to the Monte-Carlo time, i. e. number of sweeps between configurations. It is calculated in the following way: Determine an

observable on every configuration⁷, $\mathcal{O}(t_s)$, where t_s is a sweep, and assuming $\langle \mathcal{O} \rangle = 0$. Then calculate

$$\langle\langle \mathcal{O}(0)\mathcal{O}(\Delta t_s) \rangle\rangle = \frac{1}{N - \Delta t} \sum_i^{N - \Delta t} \mathcal{O}(i)\mathcal{O}(i + \Delta t),$$

where N is the total number of configurations. Alternatively, by periodicity in Monte-Carlo time, also all configurations can be included. The double bracket indicate that this average is taken over different Monte Carlo times, instead the usual ones. If the average of the operator is non-zero, it is necessary to calculate (now assuming a suitable wrapping)

$$\langle\langle \mathcal{O}(0)\mathcal{O}(\Delta t_s) \rangle\rangle = \frac{1}{N} \sum_i^{N - \Delta t} \mathcal{O}(i)\mathcal{O}(i + \Delta t) - \langle \mathcal{O} \rangle^2 \quad (4.9)$$

Usual algorithms, like the Metropolis algorithm, then exhibit usually the following behavior

$$\lim_{\Delta t \rightarrow \infty} \langle\langle \mathcal{O}(0)\mathcal{O}(\Delta t_s) \rangle\rangle \sim e^{-\Delta t_s/\tau}, \quad (4.10)$$

where τ is called the correlation time. Thus, correlation between different configurations vanish exponentially fast.

However, if $\tau \gtrsim 1$, there are still (sizable) correlations between two configurations left. There are two possibilities to deal with it. The best choice is to drop configurations, i. e. perform sweeps, but not perform measurements. Dropping N_d configurations will reduce the correlation time down to τ/N_d . Unfortunately, this may be too costly, which is especially true if fermions are involved. Then it can be shown that the statistical error including the effects of auto-correlations is larger by a factor of $\max(1, 2\tau)$ than without correlations. I. e. τ must be smaller than $1/2$, or the error increases. Thus, also N_d must be so large that τ/N_d is smaller than $1/2$, or a residual error enhancement of the same size remains.

Unfortunately, there is a number of caveats.

One is the quite obvious problem to determine τ . On the one hand, the subtraction (4.9) implies that it becomes hard to numerically extract long autocorrelation times. This can be complicated further by the fact that with a finite number of sweeps the asymptotic behavior (4.10) may not have been reached⁸, and subleading contributions can obscure the asymptotic behavior. In particular, this implies that it is never possible to obtain anything but a lower limit of τ , as long as no analytical solution for the theory exist.

⁷This already hints that it should be best a scalar observable to avoid effects due to, e. g. rotating configurations in inner or outer space.

⁸It can be shown that this is essentially dominated by the eigenvalues of the transition probability matrix, which are usually very hard to calculate.

Another one is that τ can be, and quite often is, dependent on the observable. This is intuitively clear, as an observable which depends on multiple, far-separated lattice points will only be completely changed once all lattice points have been changed often enough. An observable, which locally depends only on a single lattice site will be much quicker to decorrelate. Thus again, any determination will only give, strictly speaking, a lower estimate of τ for this observable. The best lower limit to τ is obtained when taking the largest of all τ obtained in all measurements.

A third one is that this implies that there are also correlations still with the initialization, i. e. the non-equilibrium configuration the Markov chain started from. As these are not relevant configurations any influence of them will disturb the assumption that the Boltzmann weight is the correct weight for configurations. As a consequence, usually a thermalization is performed in the beginning of a simulation. For this an initial number of configurations from N_t sweeps are dropped, and not used for measurements at all. Since thermalization is harder than decorrelation, as in this case two things have to happen at once, usually $N_t \gg N_d$ is chosen, in practice a factor 10 to 100 larger, but even larger numbers are not unheard of.

Finally, the exponential behavior (4.10) can be troublesome if the autocorrelation time becomes long. It can be shown that for an observable dominated by correlations over distances ξ then

$$\tau \sim \xi^z, \quad (4.11)$$

and thus the autocorrelation time increases as a power-law with the dynamical exponent z of the correlation length. The problem now arises if the correlation lengths become large. Especially, for light or massless particles on large volumes or close to second-order phase transitions, and thus especially close to the continuum limit, ξ becomes large. Thus, simulating closer and closer to the continuum limit becomes increasingly expensive in terms of decorrelation or statistical precision, if $z > 0$. This is known as critical slowing-down. Most standard numerical algorithms, including the Metropolis algorithm, have $z \sim 2$. It is possible to decrease this value substantially, and $z \sim 1$, or even $z \approx 0$, is possible for some systems. However, such algorithms tend then to be quite specific to the theory in question, and thus require large development, and thus human, work. It is always a careful balance between costs of additional simulation time and human effort at that point.

In addition, if the origin of critical slowing down is the approach of the continuum limit, there is an additional effect. If the physical volume should stay constant, it is necessary to increase the volume accordingly, i. e. by a factor ξ^d . Thus, the autocorrelation time for approaching the continuum limit is actually ξ^{d+z} , and thus critical slowing down can never be completely eliminated.

4.5 Systematic errors

Section 4.4 dealt entirely with problems arising from doing numerical simulations. These combine with additional errors, which are not entirely only due to the statistical evaluation, or cannot be beaten with a reasonable increase in computational power. These are called systematic errors. Despite their name, they are not systematically controlled, but merely a qualitative systematic understanding of them is possible. They are furthermore errors, which cannot be estimated based on a single simulation, but require usually multiple simulations with different lattices and lattice parameters to cope with them.

4.5.1 Lattice artifacts

Any lattice simulation is performed on a finite lattice and, due to critical slowing down, not at exactly zero lattice spacing, even if the position of the continuum limit in the quantum phase diagram is known. Examples of such systematic errors have already been discussed in various previous sections, especially 2.6.1 and 3.3. These are the dependency of the results on the lattice size and the lattice discretization: Every result from a numerical simulation will depend on both N and a . In the context of continuum theories this corresponds to a dependency on the infrared and ultraviolet regulator. These are called lattice artifacts.

Similarly to the continuum, the leading dependency on these lattice parameters can be determined perturbatively, though in this case as a series in a and $1/N$. Sections 2.6.1 and 3.3 give examples for these two types of calculations, respectively. However, a full analytical knowledge of the dependency would only be possible, if the theory would have been solved analytically, which for obvious reasons would make lattice simulations irrelevant. Thus, two problems remain: Is the volume large enough and the discretization fine enough to allow for using only a finite number of terms in these expansions? Are there (subleading) contributions not captured by a (truncated) perturbative series present which contribute in a quantitative relevant fashion?

The first question can again only be answered if a theory is exactly solved. Other than that, the best possibility is to determine results for multiple volumes and discretizations and compare the dependency with the expected dependency. If it fits, an extrapolation to infinite size and zero discretization is attempted. Even if it does not fit, because of the second question, a suitable parametrization of the dependency can be attempted and then an extrapolation can be performed. However, as is well known from mathematics, it is impossible to obtain a reliable extrapolation from a finite number of points. Thus, such extrapolations should always be considered as lower limits to the uncertainties due to systematic errors.

Of course, it is ideal if experimental results, mathematical exact statements, or at least results from other sufficiently reliable methods are available to compare to. If they fit with the extrapolations it gives confidence that the basic dependency on lattice artifacts has been understood well enough.

On the other hand, this makes life particularly annoying for new theories for which predictions have to be made. It should always be kept in mind that examples are known, where over a wide range of lattice parameters the expected qualitative asymptotic behavior was observed, but at even larger lattice sizes and finer discretizations it vanished again, only to reappear even closer to the thermodynamic limit, but with different prefactors.

A rule of thumb to estimate whether a volume is large enough is that $Nam_l \gg 1$ where m_l is the lightest mass in the spectrum. This will not work for systems with massless particles, where conclusions are notoriously hard to make. Likewise, $1/a \gg m_h$, where m_h is the mass associated with the heaviest particle, stable or not, under scrutiny. The latter is known in numerics also as the Nyquist theorem. This is particularly troublesome for new theories, where no information on m_l and m_h are yet known, forcing an iterative spiral to approach suitable lattice settings. Also, many contemporary lattice simulations in particle physics are still at a level where $2 \gg 1$ needs to be accepted on faith.

4.5.2 Overlap and reweighting

There is another problem, which arises at two points. This is encoded in (3.2) and in (4.3): Overlap. While the reasons behind this problem can be different, the conceptual consequence is quite similar.

Start out with (3.2). In any practical calculations, it will not be possible to include all possible operators. Some of these problems can be reduced, as will be discussed in section 4.6. However, there is always the possibility that the overlap $|\langle \mathcal{O}|0\rangle|^2$ of all operators with the vacuum (or any other state) is zero or, more likely, so small that it is zero within statistical errors. Then, such a state will be missed in any analysis. This is an overlap problem, as the chosen basis of operators is not large enough to statistically meaningfully cover this state. If measurements in a certain range should be performed, this problem becomes stronger the more states are in the range, as at least one operator with non-zero overlap has to be included for every channel. Since there is no a-priori way of determining operators with large overlap, the truncation error induced by the choice of basis remains an uncontrolled error.

The second overlap problem arises from (4.3). The choice of configurations is based on an importance sampling, i. e. the question which configurations have a large Boltzmann weight (4.2). This is good if the observable does not itself depend so strongly on the con-

figurations that the choice of configurations is affected by it. A trivial example where this happens is $\exp(S)$. This well-defined observable is large exactly where (4.2) is small and vice versa. Thus, trying to measure this observable with the usual choice of configurations does not work⁹. The reason is again that the support of the observable in the configurations selected during importance sampling is small compared to the total support. Due to the exponential form (4.2), this will only occur if an exponential enhancement is provided by the observable, which is strong enough to compete. This is fortunately rarely the case for the operators which are a polynomial in the fields. However, since these are also functions, it cannot a priori be excluded that such an enhancement is dynamically generated. It is almost impossible to detect such an overlap problem.

There is one situation, in which such an overlap problem is routinely relevant. While the Metropolis algorithm is suitable for all actions, it is often not very efficient. However, more specific updates can become slow if the action changes. In this case, reweighting is a possibility. For reweighting the action is split $S = S_B + S_r$. Only the part S_B is used for the generation of configurations. The contribution S_r is then included in the operator measured, i. e. $\mathcal{O} \rightarrow \mathcal{O}e^{-S_r}$. Provided $S_r \ll S_B$, this will not create an overlap problem, and thus $\langle \mathcal{O}e^{-S_r} \rangle_{S_B} / \langle e^{-S_r} \rangle_{S_B}$, where the subscript indicates the action used for importance sampling, will be the same as $\langle \mathcal{O} \rangle_S$. It makes only sense if the creation of configurations using S_B is much more efficient as the one using S , as the sampling becomes exponentially bad with an increase in S_r . Still, this happens.

4.5.3 Ergodicity

As noted in section 4.2, the Metropolis algorithm, as well as essentially all other algorithms, have a tendency to stay close to minima of the action. This behavior can be improved using overrelaxation, as described in section 4.3.2. However, this is not always sufficient. If it is not sufficient, the algorithm is not ergodic (enough), and the set of configurations may not represent the theory faithfully.

Problems with ergodicity can often be seen by monitoring some quantities as a function of Monte Carlo time. Ideally, any quantity should, after thermalization, fluctuate stochastically around the average value. If a problem with ergodicity arises, this is not the case, and the values fluctuate for a large number of sweeps around one value, and then quickly change to fluctuate around a different value. If this is the case the algorithm is not very ergodic. Of course, it may be that it is ergodic, but only very inefficient. This situation is nothing but strong residual correlations. However, determining the correlation time of section 4.4.2 would only show this if determined over multiple changes of the value around

⁹Of course, in the limit of an infinite number of configurations everything will be fine again.

which the fluctuation occurs. This may not happen quickly enough to be detectable. In fact, the jump from one value to another may happen only at too large Monte Carlo times to be visible in a simulation. Usually, the time spend at one value will increase the larger the volume, as there are more options to offset changes towards a different value by other lattice sites.

To detect such problems the best way is to start the simulation with random (so-called hot) initial conditions. If a so-far undetected problem with ergodicity arises, the different starting values will thermalize towards different values. If this happens, this is a sure sign for an ergodicity problem of the algorithm. Unfortunately, this does not work the other way around. The absence of such a signature does not imply the absence of an ergodicity problem. Fortunately, such situations are rare.

4.6 Spectroscopy

Spectroscopy is one of the primary goals of lattice calculations. Due to (3.2), it is also pretty demanding, as the error increases exponentially with time. Furthermore, it is not possible to use every operator with the same quantum numbers for a basis, and thus be able to obtain all energy levels. In practical calculations only a few, usually at most a few tens, of the infinite number of operators can be included. Fortunately, there is a number of possibilities to compensate at least part of this problem.

4.6.1 Time-slice averaging

There is a comparatively simple way to improve the statistics essentially by a factor of N_t . As particle physics theories are usually homogeneous in time, performing on every configuration

$$\langle \mathcal{O}\mathcal{O} \rangle (\Delta t) = \frac{1}{N_t} \sum_{t=0}^{N_t-1} \langle \mathcal{O}(t)\mathcal{O}((t + \Delta t) \bmod N_t) \rangle$$

creates an average over all time slices, and thereby improves the statistics. Of course, the time-slices are correlated. Thus, the measurements on individual time-slices are not statistically independent. However, the averaging does not alter the quantum numbers. Thus, the averaged operator is used as the physical operator, and thus interpreted as a single measurement on a configuration, rather than every time-slice separately.

Note that sometimes anisotropic lattices are used, possibly even with a different discretization in time and space directions, with a larger number of lattice points in time direction. This gives more to average over. Also, this can allow to have more points before

the signal-to-noise ratio becomes too bad. It therefore helps to improve the data quality as well. The downside is that this makes time special. Thus the aspect ratio of the lattice arises as an additional systematic error, as in the thermodynamic limit the lattice must become again symmetric to avoid breaking rotational symmetry. Especially, for non-scalar quantities this requires, in principle, to treat polarization directions along and orthogonal to the time direction differently, which is usually not done, despite aspect ratios of two and larger. This should be kept in mind.

4.6.2 Smearing

In most theories short distances are entirely dominated by quantum effects. Using operators which are evaluated on a single lattice site are therefore strongly afflicted by these quantum fluctuations, increasing the noise. On the other hand, often quantities living on a multiple of the lattice spacing are of general interest. Also, latest for bound states, they have usually a finite extent. Approximating them with operators, which are point-like, can also create an overlap problem.

To deal with this problem, the concept of smearing is introduced. The idea of smearing is to use extended operators rather than point-like operators. They are thus smeared over some number of lattice points. A smeared operator for the scalar field can, e. g., be constructed as

$$\phi^n(x) = \frac{1}{1 + 2(d-1)} \left(\phi^{n-1}(x) + \sum_{\mu} (\phi^{n-1}(x + \mu) + \phi^{n-1}(x - \mu)) \right), \quad (4.12)$$

where ϕ^0 is the unsmeared field. This is applied to every lattice site. Note that only nearest neighbors are involved. Therefore, smearing a single site and then continuing to the next will have the problem that the smearing is different at different positions in the lattice. Especially if n is small, this can induce systematic effects, as it signifies a preferential position. Therefore, it is best to perform smearing on a copy such that every lattice site uses the original fields for smearing. Of course, the configuration used for the next sweep in the Markov chain should not be changed, as smearing alters the Boltzmann weight of a configuration¹⁰.

By performing a number of such smearings, the operator becomes averaged over more and more lattice sites. Since the fluctuations are random and of either sign, they will average out, while the signal remains stable.

¹⁰In fact, in the limit $n \rightarrow \infty$ the result will be the vacuum, i. e. all fields vanish at every lattice site. Already much before this, the field configurations becomes classical, as all quantum fluctuations have been removed. Such configurations are actually of measure zero in the path integral, showing how smearing on the configurations of a Markov chain would distort the importance sampling.

Operators build from fields at some smearing level n are independent of operators build at a different smearing level m . Thus, using operators build from different smearing levels can also be used to create additional operators for an operator basis.

It should be noted that (4.12) is just an example of smearing. There are many other possible smearings, some particular for a given theory. They differ in the way how strongly they average. Thus, with some smearing prescription it may take (many) more smearing sweeps than with others to obtain a similar level of averaging. Also, the smearing will proceed at different speeds for different discretizations and lattice spacings. To compare the impact of different smearings or the same smearing for different lattice spacings it is best to monitor some observable quantity. If it is altered by the smearing to the same extent, i. e. has the same value, the number of smearing sweeps in both approaches are roughly equivalent. This can be understood straightforwardly: The extent of an operator smeared n times on a lattice with discretization is a . Operators with different physical extent¹¹ an , due to different discretization and different smearing, will have different overlaps with a state. Given that only a finite number of points is available, this can affect results of fits.

Besides these discrete smearing operations there exist also a continuous formulation of smearing, the so-called Wilson flow. It can be used to show that all of the smearing prescriptions are equivalent when reaching the same physical extent, as discussed above. The Wilson flow is quite similar to the concept of a continuous coarse-graining, and thus is effectively a differential equation. To apply it is essentially the numerical solution of a (partial) differential equation. This will not be detailed further here.

It is important to note that (4.12) should only be performed over spatial directions, i. e. the sum should exclude the time direction, if quantities like a temporal correlator in spectroscopy should be considered. Otherwise exactly the time dependence, which should be measured, would be washed out.

4.6.3 Variational analysis

Another possibility to improve the situation is by utilizing the insights gained in section 3.2 and 3.3. According to (3.2), every correlation function of operators is a sum of exponentials. Conversely, some linear combination of correlation functions must therefore be pure exponentials. If it would be possible to determine this basis, there would be no higher-state contaminations, and therefore already short times, where the signal-to-noise ratio is as good as possible, could be used to determine the energy-levels.

Assume for the moment that the set of operators is finite, and given by $\{\mathcal{O}_i\}$. Defining

¹¹Note that there is a usually a prefactor, which is characteristic for the smearing method.

the correlator matrix

$$\mathcal{O}_{ij}(t) = \langle \mathcal{O}_i^\dagger \mathcal{O}_j \rangle (t)$$

will give a hermitian (symmetric) matrix, if all operators are hermitian (real). For operators creating states, this can be safely assumed. Also, this requires time-slice-averaging of all operators, as otherwise different evaluation times could break the symmetry of the matrix, at least at finite statistics.

Hermitian matrices can always be diagonalized, and their eigenvalues are real. The corresponding eigenvectors will be some linear combinations of the operators themselves. In this eigenbasis, these operators have no longer any overlap with the other operators. This implies that they are necessarily pure states, and thus the eigenvalues behave like¹²

$$\lambda_i = a_i e^{-E_i t}$$

and the eigenvalues actually determine the composition of this pure base.

So far the theory. In numerical calculations, it is not possible to cover the whole (infinite) basis. Therefore, this decoupling will not be perfect, and eigenlevels will behave like

$$\lambda_i = a_i e^{-E_i t} (1 + \mathcal{O}(e^{-\Delta E_i t})),$$

where ΔE_i is the distance to the next closest energy level. Furthermore, because of the prefactors, the assignment of an eigenvalue to an energy level is not tied to the absolute size of the eigenvalue, and avoided level-crossing can occur. This can be monitored using the eigenvectors, as they should not depend on time except when such level crossing occurs.

This implies that especially for densely spaced energy regions even this so-called variational analysis will provide only an estimate. Still it is (usually) improved compared to the pure operators.

The situation improves if the coefficient of the subleading contributions can be made smaller. While this is no improvement in principle, this improves the signal-to-noise ratios at long times. It is found empirically that preconditioning in the form of the so-called generalized eigenvalue problem can do so. In this case, the eigenvalue problem is solved not for the matrix $\mathcal{O}(t)$ but for a matrix $C\mathcal{O}(t)$, where C is time-independent. Often a choice like $(\mathcal{O}_i(t_0)\mathcal{O}_j(t_0))^{-1}$ is suitable. That this indeed helps can be understood intuitively. This normalizes at a fixed time t_0 all eigenvalues to one. This is always possible, as this is just the overall normalization of the operators, which will play no role for the exponents. In the end, it corresponds to a rescaling of the operators, which is always possible. But

¹²If the size of the lattice plays still a role, this will be a cosh behavior, rather than an exponential behavior.

at this time then all exponentials are of the same size, and not hugely different. Thus, up to differences of in practical applications not more than 10 in the energy levels, all eigenvalues behave in the same way. Without this rescaling, there can be additional factors of several orders of magnitudes. Given the errors and finite precision arithmetic, this would numerically deteriorate the signal. By this normalization, this problem is avoided. Of course, to obtain the eigenvectors of the original problem then requires an appropriate rescaling of the obtained eigenvectors.

4.6.4 Resonances

Concerning resonances, everything which has been said so far for should also be applied to the determination of energy levels for counting of levels and for a Lüscher analysis. Note that smearing must be applied in such cases with great care. After all, higher energy levels oscillate quicker, and will therefore be averaged out before lower states are affected. Thus, too generous smearing could lead to a loss of signal of one or more levels, and thus of a loss of the signal of the resonances. Using therefore a mix of different smearing levels is often useful to capture all states.

Another issue is the selection of operators. It is not a-priori clear which operators have a large overlap with a resonance. Especially, if a resonance should be a collective excitation this would require operators which involve an infinite number of elementary field operators for overlap, like an exponential in the fields. Such an object would, however, likely give rise to an overlap problem. It is therefore careful balancing necessary to identify a suitable set of operators.

As a rule of thumb all field content, including those with relative momenta¹³, should be included for the ground state and all scattering states levels up to, at least, the inelastic threshold. The number of scattering levels can be judiciously reduced by reducing the physical volume. In fact, in this case a small volume can be an advantage. However, the volume should still be large enough that all other massive states are essentially at their infinite-volume value, i. e. the exponential corrections of section 3.3 should have died out already. This is usually possible to achieve.

Thus, the composition of a suitable operator basis often requires first a good understanding of ground states in the possible decay channels, if possible.

Other than that the best way ahead is a general careful analysis of the energy levels.

¹³Note that in practice discretization errors seem to play a smaller role than expected, and e. g. the use of dispersion relations correcting for the final lattice spacing are often overkill.

4.7 Phase transitions revisited

The detection, and classification, of phase transitions is essential to find the continuum limit of a theory. Performing a finite-size scaling analysis, like discussed in section 2.8, is certainly the best way. However, first and second order transitions can become quite similar, depending on the relevant critical exponent. Also, as will be discussed in chapter 7, sometimes first order transitions can be physical.

To identify first order transitions the Monte Carlo process itself can actually be used. At a first-order transition there is a coexistence of two phases. On a finite lattice, this coexistence region is actually smeared out over a range of parameters rather than just at the transition. If there are two phases, they will have differing values for observables. As a consequence, there will be configurations belonging to either phase¹⁴. In a histogram of the value of the observable in the configurations thus a double-peak structure will arise¹⁵. Observing such a structure is a signal for a first order transition.

In fact, if the Markov chain is not very ergodic, this will manifest like critical slowing down that the observable will scatter for a while around one value and then around another one. Thus, already the Monte Carlo histories can signal a first order phase transition. However, this may actually be just an ergodicity problem, and the observation of the double peak structure is much better suited. In fact, this is a signal that the physical effects of the phase transition are not sufficiently compensated for in the algorithm, and thus a bug is here turned into a feature.

Another possibility to detect a first order phase transition is by hysteresis. By tuning the control parameter in which the phase transition occurs from above the transition to below the transition during the Markov chain, the theory will stay in one phase or the other for a while even after crossing the critical value. Thus, the Monte-Carlo history as a function of the value of the control parameter will show a hysteresis loop. Again, this is using the bug of not being quickly ergodic, and an ideal algorithm would not see the effect.

However, it should be noted that an ideal algorithm would have a dynamical critical exponent of zero in (4.11). No such algorithm exists for any relevant particle physics theory. Thus, this bug can in practice always be utilized as a feature.

¹⁴On sufficiently large lattices, or sufficiently far away from the critical region, even a single configuration will have contributions from both phases, which will therefore by spatial averaging wash out the effect to be described. This needs to be offset by being close enough to the critical region at fixed lattice volume.

¹⁵Of course, this will require sufficient statistics, especially if the average value in each phases are close to each other.

4.8 Spontaneous symmetry breaking in numerical simulations

Ergodicity and a finite volume has also another consequence: Lattice simulations implement the concept of symmetry in quantum field theory quite faithfully. In particular, any non-explicitly broken symmetry is always exact, as in a finite volume and with an ergodic algorithm, if enough configurations are sampled, always all possible directions are sampled. Thus, any quantity not invariant under a symmetry transformation vanishes.

Thus, something like spontaneous symmetry breaking is not possible in a (numerical) lattice calculation. This is, however, not a bug. Rather, it really implements symmetries in quantum field theory much more literally than usually calculations. As also here spontaneous symmetry breaking, without an external limiting procedure, does not exist. At most, a theory can be metastable against any infinitesimal explicit symmetry breaking.

To obtain therefore the usual notion of spontaneous symmetry breaking requires the same prescription as in full quantum field theory: Introduce an explicit symmetry breaking, do simulations for different values of the explicit symmetry breaking, and ultimately extrapolate to zero. Only then quantities like magnetization do not vanish.

Of course, depending on the rate of ergodicity, this can be obtained with more or less numerical noise, especially at small breakings. An alternative are observables, which detect metastability, but are invariant under the symmetry. E. g., in case of the ϕ^4 model, instead of using the vacuum expectation value of an order parameter

$$\left\langle \sum_x \phi(x) \right\rangle,$$

which without explicit symmetry breaking always vanishes, it is possible to use

$$\left\langle \left(\sum_x \phi(x) \right)^2 \right\rangle, \quad (4.13)$$

which, in the infinite-volume limit, will be non-zero only in case of metastability. If there is no such metastability this quantity drops like an inverse power of volume. Thus, the situation is often straightforwardly distinguishable¹⁶.

Note that

$$\left\langle \sum_x \phi(x)^2 \right\rangle$$

¹⁶Though other lattice artifact can still fake a wrong behavior far away from the thermodynamic limit, and also a dependence like $a + b/V^c$ is possible, where $a \ll b$ leads at first to an apparently different behavior.

is not suitable, as this quantity is always non-zero, except when the field is always zero except for measure-zero parts of space-time. The latter is, however, just the vacuum.

The only reason why an explicit breaking may technically be more advantageous than using quantities like (4.13) is that operators like (4.13) generically involve a larger number of operators, and are therefore much more noisier. It then requires a careful consideration which is the better approach.

This has also implications for the search for (second) order phase transitions of section 2.8 using order parameters. These will be zero without explicit breaking, but with explicit breaking the searched-for phase transitions often degenerate into crossovers. Thus, it is either necessary to use quantities like (4.13), which are expensive, or to employ critical slowing down to make sure that the washing-out does not happen.

In the end, once more, any statement will be affected by lattice artifacts. Thus, any results will only be true up to the lattice artifacts.

Chapter 5

Gauge fields on the lattice

The formulation of gauge theories on the lattice is quite different from the one in the continuum. The main difference is that a lattice allows for a non-linear formulation, which makes the volume of gauge transformations finite, i. e. compact¹. Therefore, the path integral of a gauge theory on a finite lattice is well-defined, even without gauge-fixing.

Compactness is achieved by using instead of the gauge field A rather the parallel transporter

$$U \sim e^{iA}$$

as the elementary degree of freedom, the so-called link. Thereby, the algebra-valued gauge fields are mapped to the gauge group. Since non-Abelian Lie groups have a finite group size, in sense of the Haar measure, the links only vary over a finite range. This makes an integration over the gauge orbits without gauge fixing possible. Note that this could, in principle, also be done in the continuum, allowing for a gauge-invariant formulation. The drawback is that this induces an infinite number of tree-level diagrams, due to the exponential.

It is a far less obvious statement, and this will be discussed in detail in section 5.8, that this also allows a reformulation of the theory entirely in terms of gauge-invariant variables. However, the number of variables is infinite, and therefore this is not helpful in practice, but an important conceptual insight: Gauge fields are just auxiliary field to have a local formulation of the theory. But this is not necessary, just convenient.

¹There are also lattice formulation of gauge theories which use an explicitly non-compact version. While they have the same continuum limit, and therefore are conceptually as well-founded, as their compact counter parts, they reintroduce the problem of flat directions leading to runaway situations, and currently no efficient numerical tools are available for them in the non-Abelian case.

5.1 Abelian gauge theories

Interestingly, Abelian gauge theories on the lattice are actually more complicated than non-Abelian gauge theories. The reason is that Abelian groups are not in the same sense compact as non-Abelian groups. Because the scalar product on the generator space has zero eigenvalues, there are even for the group flat directions. Secondly, Abelian gauge theories have multiple physical charge superselection sectors, where non-Abelian gauge theories have only one. This also complicates matters. While the second problem can be tamed by suitable boundary conditions, even though these make measurements generically more complicated, the first problem has never been fully solved. The main drawback is that it is hard to find efficient ergodic algorithms. These issues become important for combined QCD+QED simulations, which start to become a major topic in lattice QCD calculations.

Here, Abelian gauge theories will only be used for conceptual purposes, and then neither of these practicalities matter.

To introduce a gauge field on the lattice, consider the following situation: Given two charges ϕ at positions x and y , how can a gauge-invariant operator be constructed? The answer is

$$\mathcal{O}(x, y) = \phi^\dagger(y) e^{ie \int_x^y dz_\mu A_\mu(z)} \phi(x) \quad (5.1)$$

where the exponential is known as Schwinger's line function. This is explicitly gauge-invariant, as can be seen from the fact that under a gauge transformation $G = \exp(ieg(x))$

$$e^{ie \int_x^y dz_\mu A_\mu(z)} \rightarrow e^{ie \int_x^y dz_\mu (A_\mu(z) + \partial_\mu g(z))} = e^{ie(g(y) - g(x) + \int_x^y dz_\mu A_\mu(z))} = G(y) e^{ie \int_x^y dz_\mu A_\mu(z)} G(x)^{-1}, \quad (5.2)$$

and since $\phi(x) \rightarrow G(x)\phi(x)$ under the same transformation.

On a finite lattice, the path in Schwinger's line function is a path along the connections between the lattice points. The shortest possibility is a straight line in a fixed direction of length a . In this case

$$e^{ie \int_x^y dz_\mu A_\mu(z)} \rightarrow e^{ie a_\mu A_\mu} = U_\mu(x) \quad (5.3)$$

where in the second step no longer a summation is performed on μ . The quantity U_μ is called a link, as it lives on the connection between two lattice sites. The link in the opposite direction is given by

$$U_{-\mu} = e^{-ie a_\mu A_\mu} = U_\mu^\dagger = U_\mu^{-1}.$$

Note that the last two equalities use the explicit fact that an element of $U(1)$ is unitary.

This implies that

$$\phi(x + e_\mu)^\dagger U_\mu \phi(x)$$

is gauge invariant. Moreover

$$\sum_{\mu} \phi(x + e_{\mu})^{\dagger} U_{\mu} \phi(x)$$

is also rotational invariant. Thus, a gauge-invariant action for gauging the ϕ^4 theory (2.23), scalar QED, is given by

$$S = \sum_x \left(\phi^{\dagger}(x) \phi(x) + \gamma (\phi^{\dagger}(x) \phi(x) - 1)^2 - \kappa \sum_{\pm\mu} \phi(x)^{\dagger} U_{\mu}(x) \phi(x + e_{\mu}) \right). \quad (5.4)$$

where translation invariance is used in the last term. Set the mass to zero and the self-coupling to zero, and thus $\kappa = 1/2d$ for simplicity. Then that this is true can be seen from using that

$$2d\phi^{\dagger}(x)\phi(x) = 2 \sum_{\mu} \phi^{\dagger}(x) U_{\mu} U_{\mu}^{-1} \phi(x).$$

Expanding then for small a

$$U_{\mu} \approx 1 + ieA_{\mu} - e^2 \sum_{\mu} A_{\mu} A_{\mu}$$

the contributions of order e^0 provide again $\phi^{\dagger} \partial^2 \phi$. Collecting everything to order e yields

$$\frac{ie}{2} A_{\mu} \phi^{\dagger} \left(\phi - \frac{1}{d} \sum_{\mu} \phi(x + \mu) \right) = ie A_{\mu} \phi^{\dagger} \partial_{\mu} \phi$$

and in the same manner with ϕ and ϕ^{\dagger} exchanged. Finally, for order e^2 , several times the same terms appear, yielding in total

$$-e^2 A_{\mu}^2 \phi^{\dagger} \phi$$

All, together, this yields

$$\phi^{\dagger} (\partial_{\mu} + ieA_{\mu})^2 \phi,$$

and thus the usual kinetic term. Hence, the (naive) continuum limit of the gauge-kinetic term is correct. Since the potential is build from gauge-invariant quantities, this is also correct in the same way.

Thus, it remains only to build the kinetic term for the gauge fields themselves, in terms of the links. This requires a gauge-invariant object, which does not involve charges. Because of (5.2), a link transforms as

$$U_{\mu} \rightarrow G(x + \mu)^{-1} U_{\mu} G(x). \quad (5.5)$$

Thus, to obtain a gauge-invariant quantity, it is necessary somehow to cancel any appearing G by a G^{-1} . The simplest quantity, which has this property is the so-called plaquette

$$U_{\mu\nu} = U_\mu(x)U_\nu(x + \mu)U_\mu(x + \nu)^{-1}U_\nu(x)^{-1}. \quad (5.6)$$

The path traced out is a square on the hypercubic lattice. In this form, it is traversed counter-clockwise, but it remains also gauge-invariant when traversed clockwise. However, by using the link definition (5.3), it can be shown that

$$U_{\mu\nu} = e^{iea^2 F_{\mu\nu}}, \quad (5.7)$$

where $F_{\mu\nu} = \partial_\mu A_\nu - \partial_\nu A_\mu$ is the ordinary field strength tensor. It is not yet a rotational-invariant object. This is obtained by summing over all possible indices which, due to the symmetry of the lattice, reduces to

$$\sum_{\nu > \mu} U_{\mu\nu},$$

Before continuing, it is convenient to rescale $A_\mu \rightarrow A_\mu/e$, removing the electric charge from the link definition, and defining $\beta = 1/e^2$. Because of (5.7), the plaquette itself would expand to the field strength tensor alone in next-to-leading order. To cancel this, the final action for Maxwell theory can be defined as

$$S_M = \beta \sum_{x, \mu < \nu} \left(1 - \frac{1}{2}(U_{\mu\nu} + U_{\mu\nu}^\dagger) \right) \approx \frac{a^4 \beta}{4} \sum_{x, \mu, \nu} F_{\mu\nu} F_{\mu\nu} + \mathcal{O}(a^5),$$

and thus the conventional gauge action. Combining this with (5.4) yields then the action for scalar QED. Defining $\beta = 1/e^2$ is not a random choice. Rather, when considering only Maxwell theory, it acts like a temperature.

The main advantage of the formulation (5.3) is that the values of the links U_μ are restricted. As an ordinary complex number on the unit circle, their absolute value is always one. This is in stark contrast to the field A_μ , which is a real number, varying from $-\infty$ to $+\infty$. This is therefore called a compact formulation of Maxwell theory or (scalar) QED. In particular, the integration variables in the path integral are not the fields A_μ , but rather the links U_μ . Thus, the integration region is finite. The integration measure is only equivalent to the usual one in the continuum limit. Otherwise, it differs by $\mathcal{O}(a)$, and thus generally the lattice theory is a different theory than the continuum one. This can be treated as an additional discretization artifact, but will, e. g., add additional $\mathcal{O}(a)$ vertices in the perturbative treatment in section 5.6.

This has an extremely important consequence when it comes to the gauge freedom of the theory. For the gauge fields A_μ , the theory is invariant under $A_\mu \rightarrow A_\mu + \partial_\mu g$, with g

an arbitrary function. Note that this function is neither necessarily continuous or needs to fall off at infinity, as is assumed in perturbation theory. In the continuum, this yields a technical problem, as there are flat directions of infinite size, the functional space of the real function g at every space-time point, where the action is constant. This yields factors of infinity. In the compact formulation, there is still a redundancy, but the gauge transformation G in (5.5) is $\exp(ig(x))$, and therefore also of absolute size 1. Therefore an integral over all possible gauge transformations for a fixed configuration, a so-called gauge orbit, is finite at every space-time point. Since the number of space-time points is finite, the total integral is finite. Thus, the path integral is finite and well-defined. Therefore the gauge-fixing, which is necessary² in the continuum to regulate the infinities, is not necessary on a finite lattice. Therefore, as long as only gauge-invariant quantities are calculated, gauge-fixing is not needed on the lattice.

However, it is possible. If gauge-dependent quantities, like the photon or (scalar) electron propagator³ should be calculated, it is even necessary as otherwise these quantities vanish or are just random (Gaussian) noise. How to do so will be discussed in section 5.5.

There is a slight complication when it comes to numerical simulations. Since different gauge copies on a given orbit give the same result for a gauge-invariant observable, it is necessary to sample every gauge orbit the same number of times, or otherwise different gauge orbits are weighted incorrectly relative to each other. In actual numerical simulations this is of little concern, as the number of configurations is so small that the chance to hit the same orbit in a Markov chain even twice can be neglected for continuous gauge groups. However, the assumption underlying the limit of an infinite number of configurations is that the Markov chain samples all gauge orbits in the same way. It also implicitly assumes that the size of every gauge orbit is the same, up to measure zero contributions. There is no reason at the current time to doubt this assumption. The situation is more involved in the case of discrete gauge theories. There, depending on the lattice size, the size of the gauge group, and the volume the possibility to hit gauge orbits more than once can become non-negligible.

There are three remarks to be made on this theory, which applies to all versions, i. e. Maxwell theory, scalar QED, and QED itself.

The Maxwell action (5.8) is up to order a^5 just the continuum action. However, starting

²Actually, it is possible even in the continuum to find formulations without explicit gauge-fixing, but so far none of them has turned out to be practically viable.

³Note that the physically observable electron or photon are not the two-point functions, i. e. propagators, of the elementary fields. They need to be dressed by a so-called Dirac phase factor, effectively adding a cloud of virtual particles. This is somewhat non-trivial, and therefore slightly outside the scope of this lecture.

from order a^5 terms like F^3 appear. Such terms induce already at tree-level additional interactions between photons. Thus, at finite a , Maxwell theory on a lattice is interacting. These interaction diminish as a goes to zero, but remain at any finite lattice spacing. In addition, there are interactions to arbitrary high orders at tree-level, suppressed only by more powers of a . Though the theory is still renormalizable when sending the regulator to infinity, i. e. a to zero, the theory is different from Maxwell theory at any finite lattice spacings. This feature will actually be generic for lattice gauge theories, because of the compact formulation. This is the price associated with using a compact formulation. Such additional interactions can therefore be considered to be additional lattice artifacts.

To remove them, it appears to be tempting to use a formulation using only the gauge fields, and thus treat them in the same way as the scalar fields and do not introduce links. This is the so-called non-compact formulation. Of course, this reintroduces the flat directions. A sufficiently small subsample, as is usual in numerical simulations, will still not hit the same gauge orbit twice. Thus, it would still be possible to do simulations without gauge-fixing. However, on any finite lattice, this is a different theory than the compact one. In particular, its phase diagram will look different. Only the continuum limit will be the same. Also, algorithms which work well for the compact formulation may not be efficient for the non-compact version, and vice versa. Thus, both theories need to be considered independently, though they are the same theory in the continuum limit.

While the two aforementioned problems apply to any gauge theory on a lattice, there is an additional problem associated with Maxwell theory and both kinds of QED. As far as it seems, the phase diagram of both the compact and non-compact versions does not have an interacting continuum limit, and the theories are trivial. While not as well established as for the case of ϕ^4 theory, there is enough circumstantial evidence that makes this result likely⁴. Thus, these theories can only be considered as effective theories, keeping a finite. But this then implies that the compact and non-compact version remain distinct. Fortunately, the impact of this distinction seems to be sufficiently small at relevant energy scales.

5.2 Non-Abelian gauge theories

Besides QED and ϕ^4 the relevant bosonic interactions of the standard model, and most theories beyond it, are non-Abelian gauge theories, so-called Yang-Mills theories.

⁴In perturbation theory, this is signaled by the appearance of the Landau pole at high energies. However, this is not an implication. The existence of a Landau pole does not necessarily imply triviality, as the case of QCD shows. Still, it is a signal which should prompt a check for triviality of a theory.

The extension of the Abelian case to the non-Abelian case on the level of the action is, as in the continuum, straightforward. In the continuum, the gauge field of Maxwell theory is upgraded to a non-Abelian gauge field by $A_\mu^a \tau^a \equiv A_\mu$, where the τ^a are the generators of any compact Lie algebra and $A_\mu^a(x)$ are d_A real fields, where d_A is the dimension of the Lie algebra, i. e. number of generators τ^a . Of course, also for the generators different representations are possible. Here, always the lowest-dimensional, usual hermitian ones will be used, i. e. for SU(2) the two-dimensional Pauli matrices and for SU(3) the three-dimensional Gell-Mann matrices.

Though the following works identically for non-(semi)simple Lie algebras only semisimple Lie algebras will be considered for ease of notation. As the generators of a Lie algebra form a vector space, this implies that at every lattice point such a vector space is attached. It is again convenient to rescale the fields by the coupling constant g , and define $\beta = 2d_F/g^2$, where d_F is the dimension of the lowest-dimensional fundamental representation. The continuum fields are then obtained as $\sqrt{\beta}A_\mu^a$ from the lattice fields A_μ^a .

As before, link variables can be defined as

$$U_\mu(x) = e^{iaA_\mu} \quad (5.8)$$

Thus, the links belong to the group⁵, rather than the algebra. But since the groups derived from compact Lie algebras are compact as well this implies, as before for the Abelian case, that the integrals over the gauge degrees of freedom are finite. However, the integrals are no longer over real or complex functions, like with the scalar or the Abelian gauge fields, but over the group. Thus, a suitable integral measure is the corresponding Haar measure.

Also gauge transformations now become group elements, i. e. G is now an element of the group as well. Then, the nominal form of the Abelian gauge transformation (5.5) remains the same. However, as the links and G are now matrices, this is a matrix multiplication. The non-homogeneous part of non-Abelian gauge transformations arises from commutators of the generators, which in contrast to the Abelian case no longer vanish. After all, formally the generator of an Abelian gauge group is just a single one, the ordinary number 1.

Likewise, to define the action it is only necessary to generalize the concept of the

⁵Note that algebras have usually more than one associated group, which usually differ by discrete factors like Z_n . Moreover, one of the groups is the universal, connected covering group, while all other groups are not necessarily connected. The choice of group can make a physical difference, especially in case of multiple gauge interactions. However, for the purpose of QCD alone, this is not the case. Since this will be the main aim in the following, just SU(3) will be assumed. The physically correct choice when embedding QCD in the standard model is actually SU(3)/Z₃, but this is less efficient to simulate numerically. In general, it needs to be specified to which group the links belong. If not stated otherwise, here always the universal covering group will be selected.

plaquette (5.6) to

$$P_{\mu\nu} = \frac{1}{2d_F} (\text{tr}U_{\mu\nu} + \text{tr}U_{\mu\nu}^\dagger) \quad (5.9)$$

$$U_{\mu\nu} = U_\mu(x)U_\nu(x + \mu)U_\mu(x + \nu)^{-1}U_\nu(x)^{-1} \quad (5.10)$$

where it is used that the representation (5.8) of the group is unitary. The normalization ensures that the value of the plaquette can range from -1 to 1 . The so-called Wilson action is then

$$S_W = \beta \sum_x \sum_{\mu < \nu} (1 - P_{\mu\nu}), \quad (5.11)$$

which reproduces the ordinary Yang-Mills action in the continuum, just like in the Abelian case. Note that the plaquette is invariant under gauge transformations, and thus the action is as well. It is also invariant under transformations from the discrete subgroup which distinguishes the groups belonging to the same algebra, in the $SU(N)$ case this is the center of the gauge group given by Z_N . This therefore establishes an additional symmetry of the theory⁶.

What has been said about additional interactions and the compact version of the Abelian gauge theory remains true for the non-Abelian one. However, Yang-Mills theory is even in the continuum limit an interacting theory, and therefore already an interesting subject in itself. Therefore, before adding scalars in section 5.11 and fermions in section 6.5, it is worthwhile to explore first this theory alone.

Note that non-Abelian Yang-Mills theories are asymptotically free, i. e. the gauge coupling at the cutoff vanishes as the cutoff is send to infinity. Since the parameters in the Lagrangian have to be interpreted as the couplings at the ultraviolet cutoff, this implies that $\beta \rightarrow \infty$ as $a \rightarrow 0$. Thus, asymptotic freedom predicts that if Yang-Mills theory is well-defined, it has to have a second-order phase transition at $\beta \rightarrow \infty$. This seems to be indeed the case, and will be discussed in more detail in section 5.7.

5.3 Gauge-invariant observables and glueballs

As in the continuum theory and the Abelian case a gauge field, or a link, are not observable quantities, as they are gauge-dependent. It is necessary to build gauge-invariant objects out of them, and thus composite operators.

For this purpose, the plaquette (5.10) delivers the blueprint: By forming a closed path of lattice sites and multiplying the links along this path creates a gauge-invariant operator

⁶Thus, formally, the symmetry group of the theory is larger than of the continuum Yang-Mills theory, and the theory is, strictly speaking, different. This does not seem to affect the continuum limit

by construction. Since the links are the only such objects available to build gauge-invariant quantities this also exhausts all possibilities. However, there are still two distinct options to consider. Since the lattice is finite, but has periodic boundary conditions, the path can either stay within the lattice volume, or wrap around once or multiple times around the lattice, by using the periodicities. Operators, which are created from closed paths inside the lattice volume are called Wilson loops. Gauge-dependent open paths are sometimes also called a Wilson line. Those derived from those traversing the boundary are called Polyakov loops, which in both cases refer to the operators or the path.

Polyakov lines, because they remain sensitive to the boundary conditions for any size of the lattice, are rather special objects. They will play a role in chapter 7, but will be disregarded for now.

The simplest operator is already the trace of the plaquette (5.9). Due to its connection to the action (5.11), the volume-averaged value of it is nothing but the internal energy of the system. By construction, it is a quantity between 0 and 2. To obtain a physical value, it still needs to be multiplied by a^{-4} , yielding an internal energy density per volume, and which is bigger than zero because of the normalization. However, because of the appearance of a in the definition of the link (5.8) the links will become closer and closer to one in the continuum limit, provided that the gauge fields do not grow like $1/a$. The latter can be proven, at least in certain gauges, most notably the Landau gauge, and is thus fulfilled. But then will the link become closer and closer to one in the continuum limit, and thus the internal energy approaches its maximum value, while a^{-4} diverges, and so does the internal energy. But this is nothing but the usual divergence of the vacuum energy, and as always, the quantity needs to be renormalized when the regulator is removed, which is just $a \rightarrow 0$ on a lattice. A renormalization scheme can always be chosen to yield zero vacuum energy.

Still, even the unrenormalized vacuum energy is a valuable piece of information. It is the most local quantity, which can be obtained. Thus, it should have the smallest autocorrelation time, and therefore the measurement of its autocorrelation time gives something of a lower limit. Furthermore, the internal energy is, of course, sensitive to the parameters of the theory. Because it is a volume average, the fluctuations of this quantity will be quite small, and thus the plaquette can be determined rather reliably, even on quite small lattices with few configurations. It is therefore a very good quantity to test, whether an implementation of an algorithm works. Of course, also the plaquette suffers lattice corrections, and therefore such comparisons should only be made on the same volumes and the same lattice spacings.

Incidentally, in two dimensions Yang-Mills theory can be solved exactly, as there it is

a theory without degrees of freedom, both in the the continuum and on the lattice. In fact, even the lattice spacing can be determined analytically as a function of β in two dimensions, and is given for an infinite volume by

$$a(\beta) = \sqrt{\frac{-\ln \frac{I_2(\beta)}{I_1(\beta)}}{\sigma}},$$

where the I_i are the i th Bessel functions and σ is a fixed dimensionful quantity, which sets the scale. This also offers a possibility to test the code. This will play a role latter in section 5.7.

Another quantity of interest is obtained from a planar, rectangular Wilson loop of size $R \times T$. As will be discussed in more detail in section 5.8, the quantity

$$V(R) = - \lim_{T \rightarrow \infty} \frac{1}{T} \ln \langle W(R, T) \rangle \quad (5.12)$$

is intimately connected with the excitation energy of a system of a static charge and a static anti-charge in the fundamental representation.

It is also possible to obtain a corresponding object in the adjoint (or any other representation). For that purpose, it is necessary to transform the links from the fundamental representation to the corresponding representation. This can be done in general, as long as distinct elements in the target representation have distinct elements in the original representation. Since the fundamental representation of the links is faithful, this is always possible. E. g. for the adjoint representation the transformation is given by

$$u_\mu^{bc} = \frac{1}{2} \text{tr} (\tau^b U_\mu^\dagger \tau^c U_\mu),$$

where u_μ is the link in the (usually not faithful) adjoint representation.

The only physical degrees of freedom in Yang-Mills theory are glueballs, i. e. bound states of gluons. Their spectroscopy can be done using the methods discussed in sections 3.2, 3.3, 3.4, and 4.6, and only the operators need to be supplied. These need to be given in terms of Wilson lines. It is now the transformation properties of the corresponding operator, and thus path, under the discrete rotation group of section 2.9.1.3, which determine the corresponding spin of the glueball operators.

For spin zero and both parities only a single representation appears, A_1 and E . For SU(2), the corresponding glueball operators are

$$\mathcal{O}_{0+}(x) = \text{tr} U_\mu(x) U_\nu(x + e_\mu) U_\mu(x + e_\nu)^\dagger U_\nu(x)^\dagger \quad (5.13)$$

$$\mathcal{O}_{0-}(x) = \sum_{\mu \neq \nu \neq \rho \neq \sigma} \text{tr} U_\mu(x) U_\nu(x + e_\mu) U_\mu^\dagger(x + e_\nu) U_\nu^\dagger(x) U_\rho(x) U_\sigma(x + e_\rho) U_\rho^\dagger(x + e_\sigma) U_\sigma^\dagger(x) \quad (5.14)$$

Thus, the 0^+ operator is just the plaquette. Note that because of the pseudo reality of $SU(2)$ there is no charge parity distinguishing states. For $SU(3)$, these operators are complex, and their real and imaginary parts correspond to positive and negative charge parity, respectively. A more involved operator is the one for 2^+

$$\mathcal{O}_{2^+}(x) = \Re\text{tr}(U_{xy}(x) + U_{yz}(x) - 2U_{xz}(x)). \quad (5.15)$$

This operator will not contain all of the multiplet elements belonging to its multiplet. It is possible to construct systematically operators of any spin, parity, and charge parity. The same is true for scattering operators. E. g. $(\mathcal{O}_{0^+}\mathcal{O}_{0^+})(x)$ represents the simplest scattering operator in the 0^+ channel. It should be noted that the 0^+ channel is in Yang-Mills theory the only channel with a vacuum expectation value, and therefore has disconnected contributions. The operators (5.13-5.15) are of particular interest in Yang-Mills theory, as the level-ordering is found to be $0^+ - 2^+ - 0^-$ for the three lowest-level, in contradiction to the naive expectation that (pseudo)scalars are lighter than tensors.

As discussed in section 4.6.2, smearing can be used to generate less noisy operators for the purpose of spectroscopy. The possible simplest smearing for $SU(2)$ is APE smearing, defined as

$$U_\mu^{(n)}(x) = \frac{1}{\sqrt{\det R_\mu^{(n)}(x)}} R_\mu^{(n)}(x) \quad (5.16)$$

$$R_\mu^{(n)}(x) = \alpha U_\mu^{(n-1)}(x) + \frac{1-\alpha}{2(d-1)} S \quad (5.17)$$

$$S = \sum_{\nu \neq \mu} (U_\nu^{(n-1)}(x + e_\mu) U_\mu^{(n-1)\dagger}(x + e_\nu) U_\nu^{(n-1)\dagger}(x) \\ + U_\nu^{(n-1)\dagger}(x + e_\mu - e_\nu) U_\mu^{(n-1)\dagger}(x - e_\nu) U_\nu^{(n-1)}(x - e_\nu)) \quad (5.18)$$

where α is a tuning parameter⁷ to optimize the effect. The quantity S , defined in (5.18), is a so-called staple. It is a plaquette in which a single link has been removed. Because of this construction, the staple transforms the same way under a gauge transformation as the removed link. This is necessary, as otherwise the smearing would create gauge-dependent results⁸. Effectively, (5.17) shows that APE smearing is again an averaging like in section 4.6.2, but this time of a link with the surrounding staple. One important restriction has to be taken into account when doing spectroscopy with smeared operators, like (5.13-5.15): Because a time dependences should be measured, a smearing which averages over

⁷A typically good choice is $\alpha = 0.55$. The extreme case $\alpha = 0$ is called cooling, and reduces quantum fluctuations fastest, but also degrades the signal very quickly.

⁸However, APE smearing is not maintaining a gauge, if one is fixed using the methods of section 5.5, but effectively performs also a staple-dependent gauge transformation.

different time-slices, including connection of time slices by a temporal link, would distort the time-dependence. Thus, for this purpose smearing should only be done with spatial hypervolumes and in spatial directions⁹.

The normalization by the determinant (5.16) is necessary, as the sum of SU(2) group elements is only proportional to a group element, and not a group element. This is called a projection into the group. This is necessary to keep within the theory. It is this last feature which makes smearing for groups other than SU(2) more difficult. In general, any averaging procedure contains summing group elements. However, for groups other than SU(2) a sum of group elements is no longer proportional to a group element. This has to be fixed by the corresponding projection. There is, however, no unique way of such a projection. A common choice is to select U_μ^n such that $\text{tr}U_\mu^{(n)}R_\mu^{(n)}$ is minimized.

Besides APE smearing many other smearing methods have been devised over time. In the end, all of them perform the same objective, albeit with different efficiencies, i. e. effect as a function of smearing progress, and have different secondary properties, e. g. whether they are analytic functions of the original links. Thus, there is no unique best choice, and secondary considerations are needed to decide what should be used for a given problem.

5.4 Numerical algorithms for gauge fields

It is possible to use the Metropolis algorithm of section 4.2 for Yang-Mills theory. A more efficient algorithm, especially with a critical slowing-down exponent of only one instead of two, is the so-called heatbath algorithm. The heatbath algorithm intends to select the new field configuration independently of the old field configuration and such that it satisfies the Boltzmann distribution. In case of a (non-Abelian) gauge theory, this is indeed possible locally, i. e. if a single lattice site is updated, the new link can be chosen such that it obeys the Boltzmann distribution for its local heatbath, formed by the surrounding links. In fact, the multi-hit Metropolis algorithm of section 4.3.1 does exactly this in the limit of an infinite number of hits. It can be shown that such an algorithm always would maintain detailed balance, but the corresponding proof will be skipped here for the sake of brevity.

In case of a gauge theory, this can actually be achieved in a single step, which is therefore highly efficient. Consider first the case of SU(2), as all other groups will be using the SU(2) procedure. This is done for each link separately, i. e., for each lattice site and each direction.

⁹If a coordinate/momentum-dependence in an additional/another direction is important, the corresponding directions should also not be smeared. Note that smearing as a site-averaging removes the high-frequency, and thus high-momentum, parts first.

The staples are the local heatbath for the link, as this is the only system with which the link is in contact. Thus, let S be the staple (5.18) for the corresponding link, and $k = \sqrt{\det S}$. This will be needed to formulate the local heat-bath.

For the next step note that $SU(2)$ is isomorphic to the three sphere, S^3 . Given any vector \vec{a} normalized to 1, the corresponding $SU(2)$ matrix m is

$$m = \begin{pmatrix} a_0 + ia_3 & a_2 + ia_1 \\ -a_2 + ia_1 & a_0 - ia_3 \end{pmatrix} \quad (5.19)$$

However, not all of this information is physically independent. After all, there is a gauge-freedom involved. Thus, it is sufficient to find a direction. Now, note that the continuum limit at $\beta \rightarrow \infty$ yields $U_\mu \rightarrow 1$, i. e. close to the north-pole to the sphere. This is corresponding to the infinite-temperature limit of the statistical system. Thus, depending on the temperature, and thus β , as formalized by the heat-bath, different distances of the direction vector to the north-pole are preferred. To accommodate it, a_0 needs to be chosen accordingly. There are different algorithms to find a good choice. Here, the so-called Creutz algorithm will be presented.

To do so, first determine a random number y by

$$y = e^{-2\beta k} + (1 - e^{-2\beta k}) r,$$

where r is a random number between 0 and 1, drawn from a flat distribution. Thus, in the limit $\beta \rightarrow \infty$ this is just r , in the limit $\beta = 0$, y is always one. A proposal for a new value of a_0 is then determined by

$$a_0 = 1 + \frac{\ln y}{\beta k}$$

and thus a_0 is essentially one at $\beta \rightarrow \infty$, but, as an expansion shows, not so at $\beta \rightarrow 0$. The rejection probability for this choice is determined as

$$P = 1 - \sqrt{1 - a_0^2},$$

and thus this proposal for a_0 is rejected with this probability, and otherwise accepted. This procedure is repeated until accepted. Note that for $\beta \rightarrow \infty$, The choice of $a_0 \approx 1$ is almost always rejected. There, it becomes hard to find a suitable result, and other algorithms, e. g. the Kennedy-Pendleton algorithm, becomes more suitable. After this, the remaining three components a_i are then selected randomly, but ensuring $\vec{a}^2 = 1$. Thus, any way of selecting a random three-dimensional (the fourth direction is a_0) direction serves this purpose. While not obvious, what is created is a new link u_μ , obtained from \vec{a} via (5.19), which satisfies the distribution

$$W \sim e^{\frac{\beta}{2} \text{tr} U_\mu S} dU,$$

where dU is the corresponding Haar measure, and thus corresponds to a local equilibrium choice. That, what makes it not obvious is the fact that there is a lot of freedom in the direction of \vec{a} , and this is explicitly accommodated here.

This completes the local update which can then be upgraded to a global update by performing it at every site, and every direction as discussed in section 4.2.

This update can be improved using overrelaxation, as discussed in section 4.3.2. An overrelaxation update is given by

$$U_\mu^{\text{new}} = \left(\frac{SU_\mu S}{\det S} \right)^{-1} \quad (5.20)$$

where S is again the staple, and U_μ is the old link. Thus, it performs a similarity transformation of the link. However, to obtain the contribution to the action, this link is again multiplied by the staple, and therefore the value of the action remains unchanged, as it must for an overrelaxation update. Still, the links are different. This is a quite efficient overrelaxation update, as it does not involve any randomness. Note that even though it may look similar to the smearing operation (5.16), it is different. There, the link is replaced by the staples, and thus, the action depends, essentially on S^2 , rather than SUS , and this changes (in fact, reduces), the action.

It now remains to extend this to groups different from $SU(2)$. This is facilitated by the fact that every Lie group has $SU(2)$ as a subgroup. The basic idea of updating groups different from $SU(2)$, the so-called Cabibbo-Marinari trick, is to update an $SU(2)$ subgroup, but always for different embeddings of this subgroup into the full group, using that there are no invariant subgroups. For this, it is necessary to rewrite every group element such that it is an element belonging to an $SU(2)$ subgroup times a remainder¹⁰. Consider, e. g., $SU(3)$. Every $SU(3)$ element can be written as a product

$$U = \begin{pmatrix} s_{11} & s_{12} & 0 \\ s_{21} & s_{22} & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} t_{11} & 0 & t_{12} \\ 0 & 1 & 0 \\ t_{21} & 0 & t_{22} \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & r_{11} & r_{12} \\ 0 & r_{21} & r_{22} \end{pmatrix} \quad (5.21)$$

The three matrices s , t , and r are matrices, which can be mapped to $SU(2)$ matrices by reunitarization, i. e. create a vector \vec{a} from, e. g., s , by setting

$$a_0 = \frac{\Re s_{11} + \Re s_{22}}{2},$$

and so on, and then create an $SU(2)$ matrix according to (5.19) from \vec{a} . These subgroups are extracted in turn for all three subgroups of the staple K . These staple submatrices

¹⁰This is the actually challenging, group-theoretical part. Usually, this requires to map the group to a complex vector space, even if the relevant representation could be real.

are then used to create a new $SU(2)$ matrix with the same heatbath algorithm as before. Note that the heat-bath algorithm does not need the previous link, and therefore this does not need to be decomposed. Call this S for the first subgroup, T for the second etc. The new link is then given by $RTSU_\mu$, where U_μ is the new link. However, to compensate for multiplying, rather than replacing, the link, requires to use for the first subgroup the original staple, for the second subgroup the staple SK , and for the third TSK . In this sense, this is a sequence of updates of the subgroup.

Rather than working in this way through the subgroups, it is alternatively possible to just select always a random subgroup. This can be done, e. g., by performing a random gauge transformation or any other means of selecting a random subgroup. The only important feature required is that this selects any subgroup of the group with the same probability. If then a sweep updates one or multiple subgroups does not matter. It is only necessary that a full update updates every subgroup in the same way, and affects every subgroup.

The same trick can be extended to any Lie group. If no explicit decomposition like (5.21) is available a choice using random transformations is always possible, as long as it guarantees to update all subgroups equally. Note, however, that the overrelaxation must also be adapted, and the simple (5.20) does not work. Again, a possibility is to overrelax all $SU(2)$ subgroups, by again selecting single subgroups in the same way, and performing (5.20) then for each.

5.5 Gauge-fixing

As noted, gauge-fixing is not necessary in numerical simulations. However, there are many reasons to determine gauge-fixed quantities nonetheless. E. g., they technically facilitate renormalization or they are of genuine interest when comparing results to continuum methods. Sometimes the behavior of gauge-fixed quantities can also be of interest in itself. Thus, there are plenty of reasons to gauge fix.

However, gauge-fixing operates quite different than in the continuum when it comes to numerical simulations. In particular, there is no Faddeev-Popov procedure or anti-field formalism. Note, however, that in analytical calculations, especially perturbation theory to be discussed in section 5.6, it works in essentially the same way as in the continuum. The only difference is that the choice of regulator, the lattice spacing, is already made. This choice is maintaining the gauge symmetry, in contrast to e. g. a cut-off regulator, and therefore does not pose any particular challenge in this context.

For numerical simulations the approach is the following. Gauge-fixing is not performed

as inserting some weight function such as to sample only part of the gauge orbits. While in principle possible no algorithm has yet been devised to do so efficiently. The reason is that the weights necessary during gauge-fixing are for relevant gauges usually such that the Boltzmann weight is modified such as to be no longer positive definite. The origin of this is that the gauge fields in principle have too many degrees of freedom, and they are essentially canceled, and thus subtracted, by the ghosts acting as negative degrees of freedom.

Thus, gauge orbits for gauge-fixed configurations are still created in the standard way. Then, gauge-fixing is performed by transforming a configuration to the desired gauge, and thereby the configuration sum (3.1) is turned into a sum over configurations satisfying a gauge condition. Since in a numerical calculations in practice never an orbit is encountered twice, not even in different gauges, there is again no problem of counting.

This leaves the problem of how to perform the gauge transformation. This is actually a more involved problem than could be anticipated. The problem is twofold. One is how to implement the usual gauge conditions. The other is the so-called Gribov-Singer ambiguity.

Start with the first one. Gauge conditions are usually given as a differential equation to be satisfied by the gauge field. An example is, e. g., Landau gauge, $\partial_\mu A_\mu^a = 0$. Other gauges are given as averages over gauge conditions. E. g. linear covariant gauges average over all possible gauge conditions of type $\partial_\mu A_\mu^a = \Lambda^a$, i. e. over all functions Λ^a , with a given weight related to the gauge parameter. The latter type is again a functional integral, and therefore would require its own Monte-Carlo process. While thus in principle solvable, this is not quite simple to implement in practice, and thus Feynman gauge on the lattice is still not really satisfactorily implemented. Therefore, here only gauges which select a single gauge copy will be considered.

It then remains to find a gauge transformation G , such that after performing (5.5), no matter whether Abelian or non-Abelian, leads to a link which associated gauge field satisfies the gauge condition, at least up to lattice corrections. Especially in the case of non-Abelian gauge theories an exact solution of the associated differential equations for matrices is usually not solvable exactly. Thus, again, numerically solutions are needed.

To give an example of such an algorithm, consider Landau gauge. Since the equation (5.5) is an equation for the link, it is useful to formulate Landau gauge using links. Consider for the following again the simplest SU(2) case. Gauge-fixing for other non-Abelian gauge groups can be constructed essentially along the same lines as with the Cabibbo-Marinari formulation for update algorithms. However, this will involve, similar to the case of smearing in section 5.3, the projection of some matrices to SU(3), and thus the same difficulties. However, the same algorithms can be used for this purpose as in the smearing

case.

The Landau gauge condition takes the form to minimize¹¹

$$E[G] = 1 - \frac{1}{2dN^d} \sum_{\mu, x} \text{tr} G(x) U_\mu G^\dagger(x + \mu) \quad (5.22)$$

over all possible G . Write $G = \exp(i\alpha\omega^a\tau^a)$, where α is an arbitrary scale parameter. Then

$$\begin{aligned} \frac{dE}{d\alpha} &= \frac{1}{2N^d\sqrt{\beta}} \sum_x \omega^a \partial_\mu A_\mu^{ag} + \mathcal{O}(a^2) \quad (5.23) \\ \partial_\mu A_\mu^a &= \sum_\mu (A_\mu(x)^a - A_\mu(x - \mu)^a) \\ A_\mu^a &= \frac{1}{2i} \text{tr}(A_\mu \tau^a) \\ A_\mu &= \frac{1}{2} (U_\mu - U_\mu^\dagger) + \mathcal{O}(a^2) \quad (5.24) \end{aligned}$$

which therefore vanishes if A_μ^a satisfies the Landau gauge condition. Note that (5.24) defines the lattice gauge field, and differs from the continuum gauge field by corrections of order $\mathcal{O}(a^2)$. Furthermore, this relation only holds for SU(2). For other groups, the traceless part of the difference has to be taken, or the trace has to be subtracted explicitly, as there the links are not hermitian.

There are many possible algorithms to find a G such that for a given original configuration U_μ the gauge-transformed links U_μ^G satisfy (5.23). Here, as an example, the simplest possibility will be considered. As noted, there is no analytic solution, and therefore an iterative method is used. As such, and working with finite precision, it will never exactly satisfy (5.23), but only approach it. It has therefore to be stopped after a sufficient fulfillment of (5.23) has been reached. Usually at least 10^{-12} is necessary¹².

The following method is called the Los-Alamos method. Start with some initial values, i. e. matrices G , at every point of the lattice. This can either be just unit matrices, or some random matrices¹³. The procedure is again organized as a change of G at every lattice site, traversing the lattice in any manner. Given some lattice site x , calculate

$$H(x) = \sum_\mu (U_\mu(x) G^\dagger(x + \mu) + U_\mu^\dagger(x - \mu) G^\dagger(x - \mu)),$$

¹¹For complex representations replace $\text{tr} U_\mu$ by $\text{tr}(U_\mu + U_\mu^\dagger)/N_c$.

¹²In fact, (5.23) is only one possibility to formulate Landau gauge, and is actually not a very precise statement, since certain contributions will not be well covered. In the literature other alternatives can be found, which are technically better suited.

¹³A random SU(2) matrix can be created by taking a random vector on S^4 , and then using the mapping (5.19).

which is essentially, the quantity with which G would be multiplied to calculate its contribution to (5.22). This matrix can be decomposed as

$$H = \sqrt{\det H} h,$$

where h is an element of $SU(2)$. Choosing a new G^n as

$$G^n = h^\dagger$$

will therefore change the contribution at x in (5.22) to the trace of the unit matrix, and therefore maximize the local decrease. However, this will change also the contribution at other lattice sites, but it can be shown that this always decreases (5.22). Thus, this is the equivalent to the heat-bath algorithm for gauge-fixing. Thus, repeating this at every lattice site will eventually lead to a fulfillment of (5.22).

Unfortunately, also such algorithms face the problem of critical slowing down. This Los-Alamos algorithm is in fact particularly bad, as it has a dynamical exponent of 2. Adding the equivalent of overrelaxation of section 4.3.2 can bring this dynamical exponent down to 1 with little effort, and essentially down to a logarithmic behavior using so-called Fourier acceleration. However, the latter has a quite large pre-factor, and therefore is only on very large lattices advantageous.

The appearance of critical-slowness indicates already a problem of (5.22): There appears to be multiple extrema. Thus, the Landau gauge condition is not unique. This is the second problem, the Gribov-Singer ambiguity. Non-Abelian gauge theories cannot be gauge-fixed using a 'local' gauge condition, i. e. gauge conditions which can be written using the fields and derivatives only. This is a generic problem, and has to do with the geometric structure of the non-Abelian gauge group¹⁴, a quite involved and not yet fully understood problem. However, the origin and details of this problem will not be the focus here, rather its implication for lattice calculations.

The problem has two aspects. The one is how to solve the problem on the lattice. This is rather straight-forward: The residual ambiguity of (5.22) has to be resolved by a prescription how to deal with the multiple minima. Many such resolutions exist, and each is a valid definition of the gauge. There are also many other gauge conditions, which are straightforward to implement on the lattice, and resolved the problem completely. The drawback is that all solutions eventually need to make some statement about that

¹⁴The non-Abelian nature manifests by the existence of the self-interactions. This has an odd consequence for Abelian gauge theories. The higher terms in (5.8) lead to interactions. They in turn make the theory 'effectively non-Abelian'. As a consequence, the Gribov-Singer ambiguity also arises in Abelian gauge theories on the lattice, but becomes irrelevant in the limit $a \rightarrow 0$.

they really do what their description proposes. E. g. one solution, the so-called minimal Landau gauge, is to average over all solutions to (5.22) with a flat weight. Proving that any algorithm devised to do so actually does so is so far impossible. On the other hand, all algorithms which can be shown to solve the problem exactly usually scale bad, often exponentially, with volume. Therefore, a proven, efficient solution is not yet there.

The second aspect is the comparison with continuum calculations, one of the motivations for gauge-fixed calculations. So far, it was not possible to construct any of the resolutions of the Gribov-Singer ambiguity on the lattice in such a way, as that it would be possible beyond doubt to implement them also in any non-perturbative continuum method. It is therefore not yet entirely clear, how to actually compare results from the lattice with continuum results.

Fortunately, all results so far strongly indicate that any resolution of this problem will only affect the very far infrared, which appears to be irrelevant for essentially all physical observables, and thus even some disagreement or unresolved subtleties in this region could very well be irrelevant in the end. Still, understanding this is essential to really fully understand (non-Abelian) gauge theories, and therefore remains open research.

A last point are the ghost fields, which are, in form of the Faddeev-Popov, or the more general anti-field, approach, introduced in the continuum to have an explicitly local formulation of the gauge-fixing. As is visible from the above, this is never done on the lattice. However, it may still be interesting to determine correlation functions of ghosts, for the same reasons as listed above. This can be done by reverting the Faddeev-Popov procedure. Originally, it was used to write the Faddeev-Popov determinant,

$$\det M^{ab} = \det \partial_\mu D_\mu^{ab}$$

as a local term in the Lagrangian. This is not done here, rather by averaging over gauge-fixed configurations in (3.1) this weight is implicitly included. Any correlation function involving ghost fields will have to have contributions with net-zero ghost number, i. e. $\bar{c}c$, as ghost number is a conserved quantum number, and therefore other correlation functions vanish. When integrating out the ghost fields, this leads to

$$\bar{c}^a(x)c^b(y) \rightarrow (M^{ab}(x,y))^{-1},$$

i. e. all ghost field pairs are replaced by inverse of the Faddeev-Popov operator¹⁵. The discretized version of this operator on the lattice reads, up to a convention-dependent

¹⁵If four or most ghost fields are involved, actually all possible pairings of the ghost have to be summed over, called in analogy to canonical quantization Wick contractions. This will actually also be true for fermions, and therefore the details of this will be postponed to section 6.8.

normalization,

$$\begin{aligned}
M(y, x)^{ab} \omega_b(x) &= \left(\sum_x (G^{ab}(x) \omega_b(x) + \sum_\mu A_\mu^{ab}(x) \omega_b(x + e_\mu) + B_\mu^{ab}(x) \omega_b(x - e_\mu)) \right) \\
G^{ab}(x) &= \sum_\mu \text{tr}(\{\tau^a, \tau^b\} (U_\mu(x) + U_\mu(x - e_\mu))) \\
A_\mu^{ab}(x) &= -2\text{tr}(\tau^a \tau^b U_\mu(x)) \\
B_\mu^{ab}(x) &= -2\text{tr}(\tau^a \tau^b U_\mu^\dagger(x - e_\mu)), \tag{5.25}
\end{aligned}$$

which is formulated as its effect on a vector. This vector can be trivially removed from the definition, but it is actually this contraction which is needed most. As the inverse of the operator has to be calculated, this needs to be done numerically. Usually, methods like conjugate gradient, or derivatives of it, are more than capable of doing so. Note that U_μ are the gauge-fixed link variables.

Note that as long as no matter fields are involved in the gauge-fixing procedure, like in the 't Hooft gauges in electroweak physics, nothing of the above is changed. Rather, only on top of determining, using (5.5), the final gauge-fixed links only also the matter fields have to be gauge-transformed using the final gauge transformation G . Since gauge transformations do not affect physics, the gauge-fixed links can also be used in the next step of the Markov chain without any consequence.

5.6 Perturbation theory revisited

Where the ghosts again appear is, of course, when doing lattice perturbation theory. Here, gauge-fixing proceeds as in the continuum. Also, as already noted in section 2.6, the Feynman rules are essentially the same, up to straightforward changes due to the lattice. However, there is one small difference for the gauge fields. As the lattice action (5.11) is formulated using the links, but perturbation theory acts upon the gauge-fields, this requires to transform from the links to the gauge-fields first, e. g. using (5.24). This difference entails slight modifications due to effects of higher orders in a .

The gluon propagator in linear covariant gauges is given by

$$D_{\mu\nu}^{ab} = \frac{\delta^{ab}}{K_\sigma^2} \left(\delta_{\mu\nu} - (1 - \xi) \frac{K_\mu K_\nu}{K_\rho^2} \right),$$

where K_μ is the improved lattice momentum (2.16) associated to the lattice momentum k_μ of (2.14) and ξ is the usual gauge parameter of linear covariant gauges. However, in contrast to the continuum case, there is also a gluon-two-point vertex,

$$\Gamma_{\mu\nu}^{ab} = \frac{(-2\pi)^4 \beta}{a^2} \delta(p) \delta_{\mu\nu} \delta^{ab}. \tag{5.26}$$

It is only contributing at zero momentum, and it diverges in the continuum limit, as there $\beta \rightarrow \infty$ and $a^2 \rightarrow 0$. In fact, it is a mass-like term. It is this term, which compensates gauge-symmetry-violations due to the regulator, and therefore is the origin of maintaining gauge-invariance with the lattice regularization¹⁶. It needs therefore to be included in lattice perturbation theory. It originates when transforming the measure from links to the gauge fields, which in the continuum limit are equivalent up to $\mathcal{O}(a^2)$ corrections. Exponentiating the latter yields, among other effects, the vertex (5.26).

The ghost propagator is comparatively harmless, as its scalar nature is not affected by the lattice, and therefore its propagator reads

$$D_G = \frac{\delta^{ab}}{K_\mu^2}.$$

Note that as in the continuum it has the wrong sign for a physical particle.

More substantial deviations arise at the three-point level. The ghost-gluon and three-gluon vertex read, respectively,

$$\Gamma_{\mu abc}^{\bar{c}cA}(p, q, k) = ig(2\pi)^4 f_{abc} K_\mu \cos \frac{ak_\mu}{2} \quad (5.27)$$

$$\begin{aligned} \Gamma_{\mu\nu\rho abc}^{AAA}(p, q, k) &= ig(2\pi)^4 f_{abc} \left(\delta_{\nu\rho} (\overline{k-q})_\mu \cos \frac{ap_\nu}{2} + \delta_{\mu\rho} (\overline{p-k})_\nu \cos \frac{aq_\rho}{2} \right. \\ &\quad \left. + \delta_{\mu\nu} (\overline{q-p})_\rho \cos \frac{ak_\mu}{2} \right) \end{aligned} \quad (5.28)$$

$$\overline{p-q} = \frac{2}{a} \sin \frac{(p_\mu - q_\mu)}{2}, \quad (5.29)$$

where momentum conservation, which holds modulo Brillouin zone translations, is suppressed, and no summation is implied. The ghost-gluon vertex (5.27) shows a lattice correction in form of a cosine-correction. For $ak_\mu \ll 1$, this correction is essentially 1. However, for momenta appreciably close to the cut-off, it is relevant. This strongly emphasizes that it would be desirable to only consider energy scales orders of magnitude smaller than the cutoff $1/a$, which is in practice usually not possible. The three-gluon vertex shows another kind of lattice correction in the form that combinations of momenta are performed on the lattice, and only afterwards the corresponding more continuum-like momenta (2.16) are obtained. Again, this creates additional lattice artifacts at scales not far away from the cutoff, i. e. wherever $\sin x \approx x$ is not a reasonable approximation.

The situation becomes even more interesting for the four-gluon vertex, which on the lattice is too lengthy to quote it here. In the continuum, this vertex is momentum-independent. This is no longer the case on a finite lattice. Because of the expansion of the

¹⁶Note that generically symmetry violations, including gauge symmetry violations, due to a regulator are not harmful when treated properly. Otherwise the violation of space-time symmetries due to the lattice would have dramatic consequences.

plaquette action (5.8), there are terms of type F^4 in the lattice action. This contains terms of type $(\partial A)^4$, which create a momentum-dependent four-point vertex. This vertex is of the same order, g^4 , as the conventional four-gluon vertex, due to the rescaling of the gauge fields with β . Because of the derivatives, it adds an (involved) momentum dependency to this vertex. Of course, this contribution vanishes as $a \rightarrow 0$, but is relevant at finite a .

In a similar way, even additional vertices are created. E. g. another four-point vertex is the two-ghost-two-gluon vertex

$$\Gamma_{\mu\nu abcd}^{\bar{c}cAA}(p, q, k, r) = \frac{g^2(2\pi)^4}{12} \{\tau^c, \tau^d\}_{ab} \delta_{\mu\nu} PQ,$$

which does not exist in the continuum theory, and vanishes in the continuum limit.

Likewise, even higher orders create further interaction vertices. In addition, it is found that individual graphs have additional divergences, and can even yield non-Lorentz-covariant terms. These are canceled between different diagrams, order-by-order. This indicates that lattice perturbation theory for gauge theories, due to the link formulation, is much more involved than continuum perturbation theory. This is the price to be paid for having a formulation convenient for numerical simulations. Fortunately, the running in QCD is fast enough that modern numerical simulations can cope with many aspects by brute force to avoid necessitating the use of perturbation theory for extrapolations or assessment of lattice corrections. For theories different than QCD this may not be the case.

5.7 Scaling and the continuum limit

An interesting result can be obtained from lattice perturbation theory for the assessment of how 'far' an actual calculation is from the continuum limit, the so-called scaling behavior. It strongly uses the fact that Yang-Mills theory (and QCD) is asymptotically free.

For this purpose, consider the running coupling. As usual, the running coupling needs to be defined within a particular renormalization scheme. For now, consider a momentum-subtraction scheme (MOM), in which it is defined to be value of the three-gluon vertex at a symmetric momentum configuration with momenta of size μ . To one-loop order, the running coupling is then given by

$$g(\mu) = g^2 \left(1 - \frac{g^2}{16\pi^2} \frac{11N_c}{3} (\ln(a^2\mu^2) + C) + \mathcal{O}(g^4) \right) (1 + \mathcal{O}(a^2\mu^2))$$

$$C = \frac{3c_3 - 2c_1}{11} \approx \frac{6\pi^2}{11N_c^2} - 9.44598,$$

and which is therefore a double-expansion in both g and a , as is always the case in lattice perturbation theory.

As always the combination $a\mu$ appears, it is now convenient to define the β -function by a derivative with respect to a , yielding

$$\beta(g, a\mu) = -a \frac{\partial g}{\partial a} = \left(-\frac{g^3}{16\pi^2} \frac{11N_c}{3} + \mathcal{O}(g^4) \right) (1 + \mathcal{O}(a^2\mu^2)).$$

Considering this as an equation for a , rather than g , and solving the differential equation, yields

$$a = e^{-\int^g \frac{dg}{\beta}} = \Lambda_L^{-1} e^{-\frac{1}{2\frac{N_c}{16\pi^2} \frac{11}{3} g^2}} (1 + \mathcal{O}(g^2)),$$

or, at next-to-leading order,

$$a = \Lambda_L^{-1} e^{-\frac{1}{2\frac{N_c}{16\pi^2} \frac{11}{3} g^2}} \left(\frac{N_c}{16\pi^2} \frac{11}{3} g^2 \right)^{-\frac{\frac{N_c^2}{(16\pi^2)^2} \frac{34}{3}}{2\left(\frac{N_c}{16\pi^2} \frac{11}{3}\right)^2}} (1 + \mathcal{O}(g^2)).$$

The newly appearing quantity Λ_L is an integration constant, and may be defined by solving these equations for it. It thus behaves like $f(g)/a$, with f the corresponding function. After fixing this quantity, e. g. by comparison to experiment, the value of a can be perturbatively determined as a function of g . Due to asymptotic freedom, given that a is fine enough, this will be its asymptotic behavior. At larger a , however, non-perturbative corrections as well as higher order corrections, will modify the dependence of a on g (and thus β). Note that $a(g)$ is a non-analytic function of g .

If $a(g)$ is found to be described sufficiently well by the above formulas, the value of β is said to be in the scaling limit, i. e. the lattice spacing effects are small enough that they are, due to asymptotic freedom, determined by perturbation theory¹⁷. Therefore, if $a(g)$ behaves perturbatively, this is usually taken as evidence for being close to the continuum limit, then also called the scaling limit. This is a genuine feature of asymptotically free theories, and, e. g., asymptotically safe theories will not show such a behavior.

More importantly, the same also holds for dimensionful quantities. Especially, for any mass m

$$m = C_m \Lambda_L$$

holds, with some coefficient C_m , which depends on the dimensionless parameters of the lattice theory, in particular later also on the bare masses of the quarks or other matter

¹⁷Note that there are some subtleties involved, as this is pure numerical evidence. It can therefore not be guaranteed that it is really under control. Also, the results will never be exactly of this form - it is just the quantitatively dominating part.

fields. Consequently, also masses will show asymptotic scaling. Measuring masses can therefore also be used to estimate how close a result is to the continuum limit. Moreover, if the Λ_L , as in the present case of pure Yang-Mills theory, does not depend on any intrinsic mass scales, this implies for the masses of two states, e. g. two glueballs, behave as

$$\frac{m_1}{m_2} = \frac{C_{m_1}}{C_{m_2}} (1 + \mathcal{O}(a^2 m^2)),$$

and thus the mass ratios of states approach constants close to the continuum limit.

5.8 The strong-coupling expansion

5.8.1 Construction

It is quite helpful to consider the equivalent of the hopping expansion of section 2.6.3 for Yang-Mills theory. Since the value of the plaquette (5.10) is bounded, the Boltzmann factor of a single plaquette approaches 1 when β approaches zero. As was discussed before, this corresponds to the limit of the coupling going to infinity, and is therefore a strong-coupling limit.

Define S_p as the contribution of a single plaquette to the action and

$$S_p = 1 + f_p,$$

then f_p vanishes as $\beta \rightarrow 0$. Since the plaquette is bounded, f_p has, in fact, to vanish uniformly¹⁸. Therefore, the full action can be written as

$$e^{-S} = \prod_p (1 + f_p) = 1 + \sum_p f_p + \sum_{p,p'} f_p f_{p'} + \dots = \sum_P \prod_{p \in P} f_p, \quad (5.30)$$

where P is any set of plaquettes. Thus, similarly to the hopping expansion, the action, and thus the path integral, can in this limit be recast as a sum over graphs with weight factors.

It is now useful to perform a so-called character expansion. For this note that the character χ_r of a matrix U^r in representation r of a group is defined as

$$\chi_r(U^r) = \text{tr} U^r,$$

and $d_r = \chi_r(1)$. Note that characters are invariant functions, as any similarity transformation of U^r will not change it. Especially, this implies that the values of a character is the

¹⁸Note that this is not true in the continuum limit of $\beta \rightarrow \infty$, as is necessary to obtain a second-order phase transition.

same for all unitarily equivalent representations. It can now be shown that any invariant function $f(U)$ of group elements can be written as

$$f(U) = \sum_r f_r \chi_r(U) \quad (5.31)$$

$$f_r = \int dU \chi_r^*(U) f(U), \quad (5.32)$$

and thus the characters form a complete basis, and r runs for simplicity over unitarily inequivalent representations.

Since the plaquette is an invariant function, it follows

$$e^{-S_p} = \sum_r d_r c_r(\beta) \chi_r(U),$$

where the c_r are coefficient functions which, in principle, can be determined using (5.32). They depend on the group. E. g. for SU(2) they become

$$a_r = \frac{c_r}{c_0} = \frac{I_{2j+1}(\beta)}{I_1(\beta)} \approx \frac{\beta^{2j}}{(2j+1)!} + \mathcal{O}(\beta^{2j+2}),$$

where the I_i are the modified Bessel functions. Reinserting this into (5.30) yields

$$e^{-S} = c_0^{6N^d} \Pi_p \left(1 + \sum_{r \neq 0} d_r a_r \chi_r(U) \right) = c_0^{6N^d} \sum_G \Pi_{p \in S_G} d_{r_p} a_{r_p} \chi_{r_p}(U),$$

where G is a graph which contains all possible combinations of plaquettes not in $r = 0$ and the support S_G of this graph is then the set of all such plaquettes. Since it can furthermore be shown that

$$\int dU \chi_r(VU) \chi_{r'}(U^{-1}W) = \delta_{rr'} \frac{\chi_r(VW)}{d_r}$$

any graph with internal boundaries can be recast into a graph which only contains an outer boundary.

In a very similar way as for the hopping expansion it is now possible to calculate any quantity of interest in a series expansion in β , as the larger the graphs the higher the power in β . However, to actually perform the calculation then requires to still calculate the integrals over the characters which stem from the path integration. In particular, the Wilson confinement criterion, which will be treated in section 5.8.2, and the fact that there is no massless excitation (i. e. glueball) is present in the spectrum can be calculated in this way.

Unfortunately, this is also already showing that the results cannot be transferred to the continuum limit. The exchange of the integration from the path integral and the

summation of the character expansion (5.31) is only possible because the functions f_p uniformly converge to zero. This is no longer true in the continuum limit, and the exchange is (in general) not possible.

There is even a physical interpretation of this process in terms of the lattice degrees of freedom. The graphs of plaquettes involved can be either (relatively) flat or can extend into all dimensions. At small β , the system is dominated by flat graphs, while at some β the importance of non-flat graphs increases. As such graphs are rougher, this is known as a roughening transition and conjectured to be of infinite order, i. e. arbitrarily smooth but still non-analytic with an essential singularity in β . This rough shape is corresponding to quantum fluctuations, as they are just random fluctuations of the fields and all associated quantities over short scales. In a sense, this can therefore be considered as the point where field-theoretical fluctuations become more important than just ordinary quantum-mechanical fluctuations.

Thus, there is no analytic connection from the strong-coupling expansion, and, strictly speaking, results in the strong-coupling expansion are meaningless for the continuum theory. Surprisingly, quantitatively, results from the strong-coupling expansion for those gauge groups, where it has been checked, are still similar to the results close to the continuum limit. However, this should not be expected to be the case.

5.8.2 Wilson criterion

To give a more detailed example of what the strong-coupling expansion can be used for, it is interesting to consider the Wilson confinement criterion.

One of the most naive ways in which to investigate confinement is to investigate the following situation, motivated by the idea of test charges in classical electrodynamics. Reduce first the problem to the quenched case, i. e. Yang-Mills theory. Place then two test-charges into the system, one with fundamental and one with anti-fundamental charge. Since test-charges can be taken to be static, this situation is completely characterized by the spatial distance between the two test charges. Connect these with a gauge field such that the total setup is gauge-invariant. Finally, measure the total energy of this arrangement as a function of the distance of the test charges.

The interesting result is that this energy has the following form,

$$V(r) = \sigma r + c + d \frac{\alpha}{r} + \mathcal{O}(\alpha^2), \quad (5.33)$$

where c and d are some constants, α is the strong coupling constant, and σ is called the string tension for reasons to become clear soon. The Coulomb-like term as well as most of the higher order corrections are what is expected in perturbation theory. In

fact, since in perturbation theory the asymptotic, non-interacting states are quarks and gluons, perturbation theory knows nothing of confinement. It is a purely non-perturbative phenomena.

The other two terms are significant. They imply that the energy is linearly rising with distance. In fact, $\sigma \sim (400 \text{ MeV})^2$ is so large that moving the two charges even the size of a proton away from each other is already very expensive, and any macroscopical scale is absurdly so. There is a restraining force associated with such a potential, which attempts to keep the charges together. This is not a necessary consequence of the requirement of gauge invariance, but a genuine feature of Yang-Mills theory.

That Yang-Mills theory can create such a behavior can be shown in the strong-coupling limit. To do so, it first requires to be a bit more explicit about the corresponding operator. It is given by the Wilson line,

$$U(C) = P \exp \left(ig \int_C ds^\mu A_\mu^a \frac{\lambda^a}{2} \right), \quad (5.34)$$

where C is a path, which starts at the first charge, moves to the second over the distance R , follows this charge for some time T , then returns to the first charge, and finally closes on itself by going back in time. It is therefore a rectangle of size RT . That the path-ordered exponential (5.34) is actually the correct expression can be seen by exponentiating the covariant derivative, which connects two fundamental quark sources over an infinitesimal distance to give a gauge-invariant expression. Its (Euclidean) expectation value

$$W = \langle U \rangle = \frac{1}{Z} \int \mathcal{D}A_\mu \text{tr} U e^{-S}$$

is gauge-invariant. On the lattice, if taking a rectangle of size the lattice spacing, it coincides with the plaquette.

From this the Wilson potential, as in (5.12), is defined as

$$V(R) = - \lim_{T \rightarrow \infty} \frac{1}{T} \log W.$$

This will ultimately yield (5.33). For this to happen, the Wilson line must behave asymptotically as $\exp(-RT\sigma)$. Thus, an asymptotic non-vanishing string tension σ implies that the exponent behaves like the area enclosed by the curve C . This is called the area law, and thus a non-vanishing σ is equivalent to confinement, and this area law is the Wilson confinement criterion. In contrast, if the exponent only scales with the length of the curve C , a so-called perimeter law, the string tension is zero, and the potential (5.33) is qualitatively the same as the one of QED, and therefore there is no confinement according to the Wilson criterion.

It is thus necessary to calculate the strong-coupling expansion of the Wilson loop of large size. Performing this calculation leads to the results that essentially a covering of the area enclosed by C by plaquettes is the dominant contribution. The contribution is then proportional to the number of plaquettes, and thus to the area. The pre-factor can then be calculated by calculating the corresponding characters. This yields for $SU(N)$ groups at leading order

$$W = \exp\left(\frac{\ln \frac{1}{3g^2}}{a^2} RT\right)$$

where a is the lattice spacing, and therefore a string tension of $1/a^2 \times \ln 3g^2$, where the lattice spacing just sets the dimensionality. Hence, the string tension is essentially given by the (large) coupling, as was to be expected given that no other parameter exist in the theory.

The drawback of this argumentation is that the strong-coupling expansion is possibly not connected analytically to the continuum limit of $g \rightarrow 0$. Hence, the proof of having an area-law in the strong-coupling limit has not necessarily any implications for the continuum theory, as reassuring as the result itself is.

Performing numerical calculations, however, show that for all practical purposes the string tension survives, provided there is no further non-analyticity involved in taking the continuum or infinite-volume limit.

5.9 Improved actions

As has been noted the action is at finite lattice spacing quite different than the continuum action, inducing even perturbatively additional terms, because of the link formulation. This is in contrast to the scalar case, and the price to be paid for having compact variables.

To reduce the impact of this problem, it is possible to work with so-called improved actions. Consider the following action

$$S = a^4 S_0 + a^6 S_1,$$

where the terms S_0 and S_1 are no longer depending explicitly on the lattice spacing a . This implies that S_1 has dimension a^{-2} . In the limit of $a \rightarrow 0$, the second term becomes irrelevant compared to the first, provided no anomalous implicit dependency of S_1 on a exists. In particular, in terms of perturbative power counting, the term S_1 is an irrelevant operator. Thus, in the continuum limit, the actions S and S_0 should describe the same theory.

Conversely, this implies that a term like S_1 can always be added to S_0 , without altering the theory in question in the continuum limit. Of course, terms with even higher powers in a can be added as well.

Now, for the Yang-Mills case, the Wilson-action (5.11) is itself a power series in a , if the links are expanded,

$$a^4 S_W = a^4 S_c + a^4 \sum_{n=1} a^{2n} S_n,$$

where $S_{n=0} = S_c$ is the continuum action. Therefore, by adding terms $-a^{2n} S_n$ to the Wilson action

$$a^4 S = a^4 S_W - a^4 \sum_{n=1} a^{2n} S_n = a^4 S_c$$

will remove the artifacts, without altering the continuum action. In practice, this is not exactly possible, as this is an infinite series. Thus, rather only the leading terms are canceled. Such actions are called improved actions. If all terms would be canceled, e. g. by finding a suitable resummation, this would be called the perfect action.

However, it is not so that improving the action is necessarily improving any observables. After all, observables contain field operators in different weights, and therefore different combinations of lower-order and higher-order operators would cancel the discretization errors. In particular, an improvement for one operator can worsen another one.

An example of such an improved action, the so-called (tree-level) Lüscher-Weisz action, can be constructed in the following way. It uses two further expressions in the action. The first is a closed Wilson loop of length six in the form of a rectangle and a three-dimensional form called the parallelogram,

$$\begin{aligned} U_{\mu\nu}^r(x) &= U_\mu(x)U_\mu(x+\mu)U_\nu(x+2\mu+\nu)U_\mu^\dagger(x+\mu+\nu)U_\mu^\dagger(x+\nu)U_\nu^\dagger(x+\nu) \\ U_{\mu\nu\rho}^p(x) &= U_\mu(x)U_\rho(x+\mu)U_\nu(x+\mu+\rho)U_\mu^\dagger(x+\nu+\rho)U_\rho^\dagger(x+\nu+\rho)U_\nu^\dagger(x+\nu). \end{aligned}$$

The resulting action, with S_W the Wilson action, is then given by

$$S = S_W + \beta \left(\sum_{x,i} \frac{5}{3N_c} \Re\text{tr}(1 - U_i^p) - \sum_{x,i} \frac{1}{12N_c} \Re\text{tr}(1 - U_i^r) \right),$$

where the sum on i includes all possible independent orientations of the corresponding Wilson loops. This action is optimized to reduce the discretization errors in certain scattering amplitudes.

Because the choice of action can be used to improve the a -dependence of observables, it makes also sense to use improved actions for theories where the action itself is already the continuum action, like in case of the ϕ^4 theory with action (2.23). However, these terms

will then deteriorate the result for the action at finite lattice spacing. Therefore, it is always important to take into account all observables which will be measured to optimize the choice of action.

It is an interesting option to combine the smearing of section 5.3 with the idea of improving the action. By smearing the operators used in the action, the system will behave better at long distances. However, care must be taken that in the limit of $a \rightarrow 0$ still the continuum action is recovered. Therefore, it is necessary that smearing needs always to be over a finite number of lattice sites, rather than over a finite physical distance, such that the operators are evaluated on a region of space-time which again shrinks to a point in the continuum limit. Still, such constructions are quite successful.

5.10 Topology

To illustrate how all the problems of lattice simulations can combine it is useful to consider a particular cumbersome observable: The so-called topological charge.

In the continuum, the topological charge is defined as

$$Q = \frac{1}{32\pi^2} \int d^d x \epsilon_{\mu\nu\rho\sigma} F_{\mu\nu}^a F_{\rho\sigma}^a,$$

where the normalization depends on the gauge-group, and is given here for SU(2). It can be shown that the topological charge is an integer in the thermodynamic limit. The reason for this being an integer is actually the global structure of space-time, i. e. how the gauge-field and action behaves at infinity. The reason is that the integrand can be recast into a total derivative, and thus Gauss law makes this a surface integral.

A lattice realization of this quantity can be obtained using the same formula together with (5.24) or, in terms of the link-variables, as a space-time sum of the operator (5.14). A quantity of similar interest is the so-called topological susceptibility, defined as

$$\chi_Q = \langle Q^2 \rangle - \langle Q \rangle^2.$$

It is connected with various statements in QCD about mass generation and in the weak interaction with baryon number violation.

Because this quantity is sensitive to infinity, it will be sensitive to the boundary on a finite space-time lattice. Moreover, it will be sensitive to the boundary conditions, therefore breaking the usual consideration that the choice of boundary conditions should no longer matter if the volume is large enough.

As a consequence, there turn up three problems in measuring this quantity on a finite lattice.

The first is that it will be some continuously distributed quantity, rather than an integer. This is an artifact because of the discrete space-time lattice and the fact that the lattice has a finite extent. Only the boundary should contribute, but the boundary is not far away, and is discrete. At infinite volume, quantum fluctuations are completely random at infinity, and would average out. This does not happen on a finite, and particularly periodic, lattice. Therefore, the averaging out of quantum fluctuations need to be done 'by hand'. This can be achieved by performing the smearing of section 5.3, as this does exactly so. However, in this case it is even necessary not only to reduce the quantum fluctuations, but to entirely remove them. Therefore, much larger numbers of smearings are necessary usually averaging many times over the full lattice. Precisely at which level to stop the smearing is a very non-trivial question, as eventually Q will become zero since eventually the smearing moves the lattice towards zero field values. This is not an entirely settled question.

The second is that because the observable is sensitive to the boundaries, and should be an integer, it is very hard to change in the Markov process. In the continuum limit, any change would need to make a big jump from integer to integer to update into a different topological sector. In a local update this becomes less and less probable the larger the number of lattice sites. This implies that the autocorrelation times strongly increase, making simulations very expensive. This problem can actually be compensated for partly by a different choice of boundary conditions, e. g. open boundary conditions.

5.11 Weak interactions and the Higgs

The simplest possibility to add matter to the gauge interactions is by coupling it to the ϕ^4 theory. This can be done for either the Abelian or non-Abelian theory in the same way. The decisive difference is again that, rather than formulating the covariant derivative using the fields, it is formulated using the links.

As a concrete example, consider the Higgs sector of the standard model. The Higgs fields ϕ are a scalar doublet¹⁹ in the fundamental representation of $SU(2)$. Thus, under a gauge transformation G , it will transform as $G(x)^\dagger \phi(x)$. Due to (5.5), the expression

$$\phi(x + \mu)^\dagger U_\mu(x) \phi(x) \tag{5.35}$$

is a gauge-invariant quantity. Since $\phi^\dagger \phi$ is so, this implies that by inserting (5.35) in (2.23)

¹⁹Actually, this case is special, as there is an additional global symmetry, the $SU(2)$ custodial symmetry. It is important to consider this additional symmetry when classifying operators, akin to flavor in QCD, but this will play no role here.

will create a gauge-invariant Lagrangian, given by

$$S = \sum_x \left(\phi^\dagger(x)\phi(x) + \gamma(\phi^\dagger(x)\phi(x) - 1)^2 - \kappa \sum_{\pm\mu} \phi(x + \mu)^\dagger U_\mu(x)\phi(x) \right. \\ \left. + \frac{\beta}{d_F} \sum_{\mu < \nu} \Re \text{tr} (1 - U(x)_{\mu\nu}) \right) \quad (5.36)$$

$$\beta = \frac{2C_F}{g^2} \quad (5.37)$$

$$a^2(2\lambda v^2) = \frac{1 - 2\gamma}{\kappa} - 2d \quad (5.38)$$

$$\frac{1}{2\lambda} = \frac{\kappa^2}{2\gamma}. \quad (5.39)$$

Expanding U_μ in a will return the usual continuum Lagrangian.

There are two caveats in working with this theory, one of conceptual and one of practical importance.

The conceptual issue arises in contact with the usual idea of a Higgs vacuum expectation value. It is possible to do perturbatively the same analysis for (5.36) as is done in the continuum, including the introduction of a Higgs vacuum expectation value. However, when doing a lattice simulations, gauge-fixing is unnecessary. Since the Higgs vacuum expectation value is actually gauge-dependent, and in some gauges always zero, this implies that a lattice simulations does not involve it. In fact, to determine it will require to fix a suitable gauge, like the 't Hooft gauge. This is just a manifestation of section 2.9: Any symmetry, which is not explicitly broken, will be manifest in a lattice simulation using ergodic algorithms. This is also true for gauge symmetries²⁰. And thus the weak gauge symmetry is unbroken. Of course, the physically observable consequences remain untouched. That the results coincide with the usual perturbative treatment is actually not trivial, but well understood. However, the explanation of this goes beyond the scope of this lecture²¹.

It is worthwhile to note that this theory is possibly trivial, an issue which may be important to provide upper limits where physics beyond the standard model may appear.

The practical issue arises with the simulation algorithm. Considering the heatbath algorithm of section 5.4, it was there necessary to determine the heat-bath in form of the

²⁰In fact, there is a proof for finite lattices, Elitzur's theorem, that it is impossible to also break a gauge symmetry spontaneously in the sense of section 2.9.

²¹The important insight is that the gauge-dependent fields, the W/Z and Higgs are not necessarily in one-to-one correspondence with the physical states, which are gauge-invariant composite objects, but for very particular reasons this is the case in the standard model for the values of the standard model parameters.

staples. Because of the coupling to the Higgs particles in (5.36) the heatbath is now a combination of the staples and the term stemming from the covariant kinetic term. Thus, the determination of the heatbath has to be adapted. In the standard model case, this is actually possible, and a modified heatbath algorithm can be constructed. However, this is due to special features of the gauge group $SU(2)$. If the theory would have a different (larger) gauge group, no heatbath algorithm is currently known, and either Metropolis-type algorithms of sections 4.2 and 5.9 are needed, or algorithms which are similar to the ones used for fermions in section 6.9.

Moreover, it is now necessary to update both the Higgs field and the gauge fields. A simultaneous change of both is usually inefficient, and thus a sweep is decomposed into usually a sweep of the gauge field followed by a sweep of the Higgs fields, or sometimes multiple sweeps for either field interlaced with multiple sweeps of the other field. It is also possible to update at every lattice point first one field and then the other. At any rate, this may degrade decorrelations, and has to be studied on a case-by-case basis.

Finally, smearing for the scalar field can no longer proceed as in section 4.6.2, in particular not like (4.12), because an expression like (4.12) is not yielding a new field which has the correct transformation properties under a gauge transformation. This can be remedied by

$$\begin{aligned} \phi^{(n)} = & \frac{1}{1 + 2(d-1)} \left(\phi^{(n-1)} \right. \\ & \left. + \sum_{\mu} \left(U_{\mu}^{(n-1)}(x) \phi^{(n-1)}(x + e_{\mu}) + U_{\mu}^{(n-1)}(x - e_{\mu}) \phi^{(n-1)}(x - e_{\mu}) \right) \right). \end{aligned} \quad (5.40)$$

The smearing for gauge fields has not to be adapted, as (5.16-5.17) maintain gauge invariance by construction, as do other smearings of gauge fields. However, in contrast to (4.12) (5.40) includes now links. Therefore, when calculating correlators smearing should not be done along time-like directions anymore.

These issues generically apply if multiple fields are involved, which now happens for the first time.

Chapter 6

Fermions

As it turns out, fermions are the bane of (numerical) lattice calculations. This is mainly due to the inherent non-local structure induced by the Pauli principle, and the fact that chiral symmetry, underlying massless fermions, is not compatible with a lattice discretization. Fermions thus require very special care. In the end, the consequence is that fermions are quite expensive when it comes to numerical simulations, and remain a substantial challenge. Although, by use of ingenious algorithms and brute force it is by now possible to simulate QCD with up to four quark flavors with physical masses.

6.1 Naive fermions and the doubling problem

The source of the problem is really the fermion itself. This can be seen by the fact that it already arises for the free fermion with Grassmann-valued fields ψ and $\bar{\psi}$. Consider its Lagrangian

$$\mathcal{L} = \bar{\psi} (\gamma_\mu \partial_\mu + m_0) \psi,$$

where the Euclidean γ_μ are hermitian¹ and are obtained from the Minkowski ones by the identification $\gamma_i = -i\gamma^i$ and $\gamma_0 = \gamma^0$.

As for the bosons in section 2.4, the fermion fields and mass are rescaled by factors of a to make them dimensionless. Using the midpoint derivative (2.9) their lattice action reads

$$\begin{aligned} S &= \sum_{nm\alpha\beta} \bar{\psi}_\alpha(n) K_{\alpha\beta}(n, m) \psi_\beta(m) \\ K_{\alpha\beta}(n, m) &= \frac{1}{2} \sum_{\mu} (\gamma_\mu)_{\alpha\beta} (\delta_{m, n+e_\mu} - \delta_{m, n-e_\mu}) + m \delta_{\alpha\beta} \delta_{nm}, \end{aligned} \tag{6.1}$$

¹This is actually a choice, and other conventions are possible.

where K is again the inverse propagator.

Using the rules for Grassmann integration, the propagator in momentum space can be calculated and yields

$$\begin{aligned} \langle \bar{\psi}_\alpha \psi_\beta \rangle(p) &= K_{\alpha\beta}^{-1}(p) = \frac{-i(\gamma_\mu)_{\alpha\beta} P_\mu + m\delta_{\alpha\beta}}{P_\mu^2 + m^2} \\ P_\mu &= \frac{1}{a} \sin(p_\mu a) = \frac{1}{a} \sin\left(\frac{\pi n_\mu}{N_\mu}\right). \end{aligned} \quad (6.2)$$

The important difference to (2.15) is the differing factors of 2 in the definition of the momentum. They arise from the first-order nature of the fermion action, and are the source of the differences to the scalar case.

The consequences are severe. In the scalar case the expansion of K_μ for small a , which is the equivalent object to P_μ in the fermion case, in (2.16), yielded only “small” (compared to $1/a$) contributions for small lattice momenta, $n\pi/N \ll 1$. Here, this also appears for $n\pi/N \sim \pi$, because of the factor 2. Thus, the propagator (6.2) has additional poles. There are 2^d to be precise, one for every combination of where the momenta are either close to zero or close to the edge of the Brillouin zone. Thus, despite the action seems to go to the correct continuum limit, it describes far too many fermions.

The reason for this is that a quantum theory is not only defined by the action, and therefore statements based just on the action are insufficient in the quantum case. Also the measure of the path integral contributes. And here the origin of this phenomenon originates. Technically, the additional poles occur because at least one momentum is of order the cutoff. Therefore, there is always a light particle at energies at the cutoff, and therefore probes the existence of a cutoff, and therefore the lattice structure. Therefore, the regulator always play a role. Why is this so?

In the continuum theory, the situation is somewhat different. There, when sending the regulator to infinity, renormalization is necessary. The renormalization, however, induces a so-called anomaly, i. e. it breaks one of the symmetries of the theory. For fermions, this symmetry is the so-called (global) axial U(1) symmetry, i. e. invariance of the action under

$$\psi \rightarrow e^{i\gamma_5\theta}\psi,$$

where θ is an arbitrary parameter, and $\gamma_5 = \gamma_0\gamma_1\gamma_2\gamma_3$ is the Euclidean version of the γ_5 matrix. For any finite regulator, the theory is actually not breaking the symmetry. The appearance of the doublers can be traced back to this. A single fermion pole at the origin would actually also violate the symmetry, as this kind of rotation transforms the 2^d poles into each other. Thus, the other poles must be present to keep the symmetry realized at

finite cutoff. It is thus a problem of the exchange of limits². This already hints at the resolution of the problem: It appears necessary to somehow break the symmetry already at finite cutoff.

To understand the problem, it is best to trace out the steps leading to the appearance of the factor 2. Going backwards, it turns out that it originates from using the midpoint derivative, (2.9), which already shows the explicit factor of 2. The reason for choosing the midpoint derivative was that it is manifestly (anti)hermitian, while left-derivatives and right-derivatives transform into each other by performing a hermitian conjugation, therefore not yielding a hermitian Lagrangian. This problem did not appear for scalars, which are described by second-order derivatives, of which the discretization is both symmetric and hermitian. The first-order nature is actually imposed by the fact that the fermions require a first-order description. All of this points to a very fundamental problem of fermions.

This is indeed the case, and formulated in the Nielsen-Ninomyia theorem, which proof is beyond the scope of this lecture: It is impossible to have at the same time, in a lattice formulation of fermions, chiral symmetry, locality, hermiticity and translational invariance.

Thus, it is necessary to give up at least one of those to make progress. Studying the list, chiral symmetry appears to be by far the most acceptable sacrifice. This is indeed what is done in most modern formulations of fermions on the lattice, though all other options have been studied as well. In the end, it is only the continuum limit, which is important. And in some cases, it was possible to work with a violation of one of the other properties, and still recovering a decent continuum theory. But the, by far, simplest, most general, and most successful approach is to loose chiral symmetry on a finite lattice.

6.2 Wilson fermions

There is a large number of ways how this can be done. Of these, only three will be presented here. One are the Wilson fermions to be discussed in this section. In a sense it is the most brutal version to deal with the problem, but this is balanced by the fact that it is the cheapest one in terms of computation time. In the section 6.4 a solution will be shown, which is arguably the most elegant one, albeit at much larger computational cost. Both versions together cover more than half of all contemporary lattice simulations.

One of the conditions remaining is that the action should still be the same in the continuum limit. Thus it is a good strategy to modify the action by higher-dimensional,

²Incidentally, the anomaly has observational consequences in QCD, and this could be taken as an indication that a continuum approach is necessary to describe physics, at least without (quantum) gravity.

so-called irrelevant, operators, which have positive mass dimension, and therefore vanish in the continuum limit.

For Wilson fermions, this amounts to change the action to

$$S = S_N - \frac{r}{2} \sum_n \bar{\psi}(n) \partial_\mu^f \partial_\mu^b \psi(n)$$

where S_N is the naive fermion action (6.1), and r is a new parameter, the so-called Wilson parameter. This is a term with mass-dimension five. For the theory to be renormalizable requires therefore that the term is taken to be proportional to a , and therefore to vanish in the continuum limit. Thus, r has to be chosen such that this is achieved.

To see how this comes about, it is useful to determine again the propagator. As the action is still quadratic in the fermion fields, this can be done explicitly, and yields

$$K_{\alpha\beta}(n, m) = (m + 4r) \delta_{nm} \delta_{\alpha\beta} - \frac{1}{2} \sum_\mu ((r - \gamma_\mu)_{\alpha\beta} \delta_{m, n+e_\mu} + (r + \gamma_\mu)_{\alpha\beta} \delta_{m, n-e_\mu}). \quad (6.3)$$

Note that the matrix-version of r is just the unit matrix multiplied by r , but forming this combination is convenient. As r also appears in the mass term, the mass term no longer vanishes at $m = 0$. Therefore, the additional term breaks explicitly chiral symmetry even for massless fermions, as was required.

To interpret the consequences, it is useful to study the propagator in momentum space, which reads

$$\langle \bar{\psi}_\alpha \psi_\beta \rangle(p) = K_{\alpha\beta}^{-1}(p) = \frac{-i(\gamma_\mu)_{\alpha\beta} P_\mu + M(p) \delta_{\alpha\beta}}{P_\mu^2 + M(p)^2},$$

which looks very much like (6.2), except for the appearance of a momentum-dependent mass

$$M(p) = m + \frac{2r}{a} \sum_\mu \sin^2 \frac{ap_\mu}{2},$$

where the lattice parameter has been reinstated for convenience. This mass shows distinctively different properties inside the Brillouin zone and at its boundaries. Since sufficiently far inside the sine expands to a , the second term vanishes for $a \rightarrow 0$. However, at the boundaries, the sine approaches 1, and therefore the mass diverges. Thus, the doublers are made infinitely heavy and are therefore removed from the spectrum.

Note that at any finite value of a the actual mass is r dependent, but can be made zero at $r = -m/4$. Once interactions are included, this is no longer analytically possible, and the actual value required for massless fermions has to be determined numerically, and depends on a .

6.3 Staggered fermions

The search for useful solutions to the problem of chiral symmetry on the lattice has led to a large number of fermion formulations. As has been explicitly demonstrated for Wilson fermions, all of them have conceptual drawbacks until arriving at a conceptual clean solution to be discussed in the next section. However, a conceptual clean solution is expensive in terms of computing time. Thus, many alternative formulations, which are compromises, have been proposed and used. Most of them have never achieved widespread use, except for Wilson fermions and the so-called staggered fermions.

The latter are technically much more involved than the Wilson fermions, but restore part of the chiral symmetry. This is achieved by having $2^{d/2}$ degenerate fermions for a single original fermion, i. e. four in four dimensions. This implies the existence of a corresponding flavor symmetry. To separate this symmetry from the actual flavor symmetry, it is usually called taste symmetry. By placing the fermions on a sublattice of half the lattice spacing, with one flavor per lattice site only, the ensuing theory has a diagonal subgroup of the original flavor and chiral symmetry, which is conserved in the lattice approximation. Especially, this symmetry can be spontaneously broken, leading to the desired excitation spectrum.

The downside of this formulation is that only theories having fermions in sets of four degenerate ones can be studied. QCD is not even approximately of this type, as the charm quark is heavy.

As an ad-hoc compensation of this problem, rather than using the quark propagator, its fourth root is used, the so-called rooting trick. Though this works on a finite lattice without gauge interaction, there are conceptual problems with gauge interactions and in the continuum limit. But quantitatively these conceptual problems seem not to matter for essentially all quantities of interest, and therefore this approach has found widespread use, especially as it is computationally quite cheap.

Still, in the long-run, with more and more readily available computing power, these will drop out of use in favor of Ginsparg-Wilson fermions to be discussed next. In addition, because of the conceptual problems, using staggered fermions with the rooting trick, remains a debated issue. This led to a renaissance of the use of Wilson fermions in recent times.

6.4 Ginsparg-Wilson fermions

The Nielsen-Ninomyia theorem forbids a real solution of the problem. However, it turns out that there is an alternative way of approaching the problem of chiral symmetry on the lattice. So far, all attempts have been focused on making the breaking of chiral symmetry by the lattice less dramatic for small a . The alternative approach is to not insist on having chiral symmetry on the lattice, but rather having a symmetry on the lattice, which becomes the original chiral symmetry in the continuum limit. In a sense, the problem is evaded, rather than solved. This should not be considered a cheap exit. After all, the Nielsen-Ninomyia theorem shows that it is impossible to solve the original problem.

The advantage of this route that there is always a symmetry, and thus spontaneous symmetry breaking is possible even for a finite lattice. The downside will be that this approach is that it is substantially more expensive than, e. g., Wilson fermions. This is the reason why not all calculations nowadays are done using this formulation. However, whether the original symmetry is explicitly broken or replaced are on any finite lattice conceptually equally good³.

The starting point for the Ginsparg-Wilson construction is the fermion action

$$S = \sum_{xy} \bar{\psi}(x)(D(x,y) + m\delta_{xy})\psi(y)$$

where the Dirac operator is now required to have two properties. The first is that for $a \rightarrow 0$ it reproduces the continuum Dirac operator, to generate the correct continuum limit of the theory. This is already true for the Wilson Dirac operator, and even the naive Dirac operator.

The second requirement is the non-trivial one⁴

$$\gamma_5 D(x,y) + D(x,y)\gamma_5 = a \sum_z D(x,z)\gamma_5 D(z,y). \quad (6.4)$$

If the right-hand-side would be zero, this would be just the ordinary requirement $\{\gamma_5, D\} = 0$ implying continuum chiral symmetry. Since D should reproduce the ordinary Dirac operator in the continuum it is not singular in $1/a$, and therefore the right-hand side vanishes as $a \rightarrow 0$, recovering chiral symmetry in the continuum limit. So far, so good. However, several questions need now to be answered.

³There are some subtleties in the continuum limit if the chiral limit is considered, i. e. zero fermion (quark) mass at the Lagrangian levels. These arise because of analyticity properties of the limit. However, these play no role if the symmetry is already explicitly broken by finite quark masses, and therefore irrelevant, e. g., in QCD.

⁴In general, γ_5 can be replaced $\gamma_5 R$, where R is any local, non-singular matrix, proportional to the Dirac-1, though it may differ from one in other spaces, e. g. in flavor space. This will not be used here.

The first is whether the theory is still well-defined, i. e. local. Rewriting this as a matrix equation also in space-time indices, the relation (6.4) can also be written as

$$\{\gamma_5, D^{-1}\} = a\gamma_5.$$

Since the right-hand side is proportional to a unit matrix in space-time indices, the breaking is even ultralocal, answering this question affirmatively.

The second is, whether there are any solution at all. This is also true, as

$$D = \frac{1}{a}(1 - V)$$

with

$$\begin{aligned} V^\dagger V &= 1 \\ V^\dagger &= \gamma_5 V \gamma_5 \end{aligned}$$

solves (6.4). These are rather weak conditions, and it is straightforward to find solutions.

The final question is actually less trivial to answer: Does this Dirac operator, called the Ginsparg-Wilson-Dirac operator, have doublers. This is in general the case. However, there are solutions, which have not. A particular example is the Neuberger-Ginsparg-Wilson-Dirac operator,

$$\begin{aligned} V &= A(A^\dagger A)^{-\frac{1}{2}} \\ A &= 1 - aK^W \end{aligned}$$

where K^W is the Wilson-Dirac operator (6.3). This is stated here without proof. The appearance of the square-root of a matrix makes this a very involved numerical problem, which is much more expensive in terms of computing time than Wilson or staggered fermions. In section (6.5), when QCD is introduced, the Dirac operator needs only to be replaced by one with minimal coupling to continue working for QCD.

It is useful to study the features of this symmetry further. The actual symmetry implemented is

$$\begin{aligned} \psi &\rightarrow e^{i\theta\gamma_5(1-\frac{a}{2}D)}\psi \\ \bar{\psi} &\rightarrow \bar{\psi}e^{i\theta\gamma_5(1-\frac{a}{2}D)}, \end{aligned}$$

i. e. it differs from an ordinary chiral transformation by the term $i\gamma_5 aD/2$. This deviation vanishes again as $a \rightarrow 0$.

Another advantage of this formulation is coming from a quite different angle. In the standard model chiral symmetry is actually quite different than in QCD. Because the

weak interactions only affect left-handed fermions, chiral transformations mix in a non-trivial way with weak gauge transformations. If chiral symmetry would be explicitly broken by the lattice regulator, this will propagate to the gauge symmetry, leading to an anomalous breaking of the weak gauge symmetry. Since a gauge anomaly renders a theory physically useless, observables depend on the choice of gauge when quantizing such a theory, this makes a lattice description of the weak gauge interaction with Wilson and staggered fermions impossible.

Having such a replacement symmetry with the correct continuum behavior available appears at first to solve this problem. Unfortunately, this is not the case. Though one can write down some associated gauge symmetry, it has not the correct continuum limit to act as a discretization of the standard model. So far, a solution has not been found, and this question is subject of current research. It seems possible that a unbroken alternative symmetry, which becomes the weak gauge symmetry in the continuum limit, exists.

6.5 QCD

After having fermions available, the next step is to study the corresponding gauge theories. Here, this will be restricted to QCD-type theories, i. e. coupling of (mass-degenerate) fundamental fermions to a non-Abelian gauge theory, usually of $SU(N)$ type. More complicated theories are straightforward extensions, as long as the fermions are coupled vectorially to the gauge fields. Chiral couplings, like in the weak interactions, are more problematic. On the other hand Abelian gauge theories, like QED, are a simplification⁵.

Turning free fermions into QCD follows essentially the same line as for the scalar fields in section 5.11, i. e. by promoting the fermions to color vectors in an appropriate representation, and introducing a link into the derivative term of the fermion action.

For simplicity, here only one class of fermions, the Wilson fermions, will be studied in detail. The full action and corresponding discussion for essentially every type of fermions can be found in the literature.

The resulting action is

$$S = S_{\text{YM}} + \sum_{f,x} \left(\bar{\psi}_f(x) \psi_f(x) - \kappa_f \sum_{\mu=\pm 1}^{\pm 4} \bar{\psi}_f(x + \mu) (r_f + \gamma_\mu) U_\mu(x) \psi_f(x) \right)$$

where S_{YM} is the Yang-Mills action and f are the flavors, which can have different masses and therefore different hopping parameters and Wilson parameters.

⁵Note that existing results on QED alone seem to suggest that it is a trivial theory in the sense of section 2.7.

In the same way as for the scalar particle an expansion in a of the links yields the covariant derivative. Also, in the same way this implies the existence of an involved tree-level structure in terms of the gauge fields at finite a .

Especially for numerics it is useful to note that the quarks only appear as a bilinear. Since they are Grassmann fields, it is possible to perform the functional integral explicitly, yielding⁶

$$\int \mathcal{D}\bar{\psi}\psi A e^{-S} = \int \mathcal{D}A \prod_f \det D_f e^{-S_{\text{YM}}}. \quad (6.5)$$

Herein the determinant of the Dirac operator D for every flavor f , $\det D_f$, appears. In particular, if the flavors are degenerate, the product is $(\det D)^{N_f}$.

This result is crucial for the development of fermion algorithms in section 6.9, as Grassmann numbers are not representable in a computer, while the fermion determinant is made of ordinary numbers⁷, and thus can be.

A second use is that it allows for a, at first sight, quite drastic approximation. Setting $\det D = 1$ gives the quenched approximation. Thus, evaluating operators involving fermions fields, see section 6.7, uses gauge fields weighted by the Yang-Mills action. Since algorithms to create configurations including the determinant are much more demanding in terms of computing time, this is a drastic simplification. While there is no justification for such an ad-hoc approximation possible, it turns out to be not a bad one. In fact, many quantities, especially static ground-state quantities, are within 10% of their values in actual QCD. This gets even better if the quarks are heavier than in ordinary QCD. This already explains why the approximation is working so well based on physics. Because of the mass of the quarks they are not relevant for long-range structure. However, static bound-state properties are long-range effects, and are therefore only weakly affected. Due to dynamical chiral symmetry breaking in QCD this effect is even amplified, as it provides every quark with an additional effective contribution to its mass, making the approximation even better.

⁶Note that there are subtleties involved with zero eigenvalues of the Dirac operator, which will not be detailed here.

⁷Note that this has nothing to do with QCD, just with the possibility of integrating out the fermions. If the theory should have terms in the action of higher power than two in the fermion fields it is possible to use a (generalized) Hubbard-Stratonovich transformation to introduce auxiliary fields to reduce the polynomials to bilinears, and therefore perform the same exact integration.

6.6 Perturbation theory for fermions

Just as with the scalars in section 2.6, there are now possible ways to treat QCD (or fermions in general) perturbatively. One is an ordinary expansion in the coupling constant, while the second one is also a hopping expansion, or equivalently a heavy-mass expansion.

Consider first the conventional perturbative expansion. This requires to supplement the Feynman rules of section 5.6 by the quark propagator and the quark-gluon vertex.

Except for the color structure, the quark propagator is essentially the one of the free quark. Thus, its propagator reads

$$D_{\psi\psi\alpha\beta}^{ab} = \delta_{ab} \left(\frac{1}{i \sum_{\mu} \gamma_{\mu} \sin p_{\mu} + M_f(p)} \right)_{\alpha\beta}$$

$$M_f(p) = m_f + 2r_f \sum_{\mu} \sin^2 \frac{p_{\mu}}{2},$$

where again the (flavor-dependent) Wilson term appears. Thus, in comparison to the continuum case, already at tree-level the mass function is momentum-dependent, and quite strongly so. Thus, in comparison to the rather mild lattice correction of the gluon and ghost in section 2.6 the fermions are much more sensitive to the underlying lattice.

Even more interesting is the case for the quark-gluon vertex,

$$\Gamma_{\alpha\beta\mu}^{\bar{\psi}\psi A f f' i j a}(p, q, k) = -ig(2\pi)^4 \delta_{ff'} \tau_{ij}^a \left(\gamma_{\mu} \cos \frac{(p+q)_{\mu}}{2} - ir_f \sin \frac{(p+q)_{\mu}}{2} \right)_{\alpha\beta}.$$

It has not only a further (real) part, but actually an additional different Lorentz structure, which arises from the Wilson term. In fact, since it is scalar, it is a term which explicitly breaks ordinary chiral symmetry. Because of its momentum structure, this will vanish again in the continuum limit, recovering the correct (perturbative) chiral symmetry of the interaction. Nonetheless, this shows how the explicit breaking of chiral symmetry manifests due to the lattice regularization not only at the level of creating an additional mass, but in the very structure of the interaction.

Other than that, ordinary perturbation theory proceeds with fermions as without fermions.

Consider in the following N_f degenerate quarks. The second type of expansion is the hopping expansion. For this, write the Dirac operator as

$$Q(x, y) = \delta_{xy} - \kappa \sum_{\mu} \delta_{x, y+\mu} (r + \gamma_{\mu}) U_{\mu}(x) = \delta_{xy} - \kappa M(x, y)$$

where M is called hopping matrix. Rewrite the fermion determinant as

$$\det D = e^{\text{tr} \ln(1 - \kappa M)}.$$

This logarithm can be expanded as

$$\mathrm{tr} \ln(1 - \kappa M) = \sum_l \frac{\kappa^l}{l} \mathrm{tr} M^l.$$

Since the hopping matrix has a simple structure, the powers can be calculated explicitly as

$$\mathrm{tr} M^l = N_f \sum_{x_i \mu_i, i=1 \dots l} \delta_{x_1, x_l + \mu_l} \delta_{x_l, x_{l-1} + \mu_{l-1}} \dots \delta_{x_2, x_1 + \mu_1} \mathrm{tr} (U_{\mu_l}(x_l) \dots U_{\mu_1}(x_1)) \mathrm{tr} ((r + \gamma_{\mu_l}) \dots (r + \gamma_{\mu_1})) \quad (6.6)$$

Because of the δ s, this is always a closed Wilson loop, and therefore on a hypercubic lattice only even l contribute, weighted with a spinor trace. E. g., up to order 4 this gives

$$\begin{aligned} \mathrm{tr} \ln(1 - \kappa M) &= N_f V (48\kappa^2(r^2 - 1) \\ &\quad + \kappa^4 \left((r^2 - 1)^2 + 8(r^4 - 2r^2 - 1) \frac{1}{V} \sum_{\text{Plaquette}} \Re \mathrm{tr} U_{\text{Plaquette}} \right) + \mathcal{O}(\kappa^6)). \end{aligned}$$

Thus, this shifts the gauge action by an offset, therefore effectively shifting β . Thus, the presence of the fermions act as changing the Boltzmann factor such that different configurations are weighted differently.

This should not come as a surprise. After all, the integration range of the path integral of the gluon field did not change when introducing the quarks. The only thing which can change is how important the gauge fields are. This is true in general - the presence of matter just weight the gauge field configurations differently, but does not change the integration range in the path integral⁸

6.7 Fermionic expectation values

Expectation values for fermions are⁹, conceptually, just the same as for bosons: Expectation values of products of the fields. While in an analytical calculation this is all there is to it, this is quite different when performing numerical evaluations as, as noted, Grassmann

⁸While in QCD the presence of the quarks is a second-order effect, as can be seen from the success of the quenched approximation, this is quite different in case of adding the Higgs in section 5.11: If the Brout-Englert-Higgs effect is active, the gauge bosons act as having a mass. The gauge fields describing this behavior are already included in the pure Yang-Mills case, but their effect cancels out, given the (pseudo)massless gluons of Yang-Mills theory. It is the skewing by the Higgs field in the sampling procedure which brings the other behavior to be dominating.

⁹Note that essentially everything, except for the Lorentz structure, said in the following also applies to determining expectation values for the Faddeev-Popov ghost of section 5.5.

numbers cannot be represented in a computer. However, the tricks necessary to do so numerically are sometimes also useful analytically, which is why they are discussed before section 6.9.

The most important insight is that the relevant expectation values are always conserving fermion number¹⁰. Thus, the number of fermion and anti-fermion operators must always match in an expectation value. Since the action is a fermion bilinear¹¹, the fermionic integration can always be performed exactly, as they are of Grassmann-Gaussian type. Doing so yields therefore expressions involving the Dirac operator only.

Performing the explicit calculation is essentially nothing but performing Wick contractions in canonical quantization, and this will therefore not be detailed here. For simplicity, here the case of N_f degenerate flavors will be considered. The final results is

$$\begin{aligned} & \langle \psi(y_1)\bar{\psi}(x_1)\dots\psi(y_n)\bar{\psi}(x_n)\mathcal{O}(A) \rangle \\ &= \int \mathcal{D}\bar{\psi}\psi A\psi(y_1)\bar{\psi}(x_1)\dots\psi(y_n)\bar{\psi}(x_n)\mathcal{O}(A)e^{-S_{\text{YM}}-\sum\bar{\psi}Q\psi} \\ &= \int \mathcal{D}Ae^{-S_{\text{YM}}}\mathcal{O}(A)(\det Q)^{N_f}\sum_{z_1\dots z_n}\epsilon_{y_1\dots y_n}^{z_1\dots z_n}Q^{-1}(z_1,x_1)\dots Q^{-1}(z_n,x_n) \end{aligned} \quad (6.7)$$

where \mathcal{O} is any operator of the gluon fields only, any gauge-index contractions have been suppressed, and a canonical order of the quark fields has been introduced, which can always be achieved by transpositions (and corresponding powers of -1). Q is given in¹² (6.6) and

$$\epsilon_{i_1\dots i_n}^{j_1\dots j_n} = \begin{cases} 1 & \text{even} \\ -1 & \text{if } (j_1, \dots, j_n) \text{ is an odd permutation of } (i_1, \dots, i_n). \\ 0 & \text{no} \end{cases}$$

If the flavors are not degenerate, Q becomes also a matrix in flavor space. Without flavor-changing interactions, it is, however, block-diagonal, and therefore the determinant factorizes, and also the inversions become simpler.

The main problem is now how to calculate both the effects of the determinant¹³ and the inverse of the Dirac operator. The first problem will be relegated to section 6.8, as it is relevant for the generation of the configurations.

The answer to the second is both straightforward and involved at the same time. In principle, it just needs to require to invert a matrix numerically. However, the matrix is

¹⁰There exist theories without this feature. They will not be considered here, and would need quite different approaches.

¹¹See footnote 7 for other theories.

¹²For the ghosts, just replace Q by the Faddeev-Popov operator (5.25).

¹³Which is simply set to one in the quenched case, but leaving the sum untouched to obtain quenched expectation values for the fermions.

$V \times d_F \times N_f \times d_D$ -dimensional, where d_D is the dimension of the representation of the Clifford algebra for the fermions. For Dirac spinors in four dimensions $d_D = 4$. It is the factor V which makes this inversion numerically involved. The one advantage is that Q , as can be seen from (6.6), is a sparse matrix, inheriting the (anti-)hermiticity properties of the Euclidean Dirac matrices. Conjugate gradient-type algorithms are therefore suitable to provide a good estimate.

There are two technical problems. One is minor, and is the choice of stopping criterion. Though conjugate gradient algorithms with finite precision arithmetics are subtle when it comes to convergence properties, the Dirac operator is quite well-behaved so that already rather simple stopping criteria are sufficient.

The second problem is that a conjugate-gradient type algorithm operates by applying iteratively the matrix to a vector to construct the action of the inverse matrix on this vector. This result is then applied in an expression like (6.7) to other quantities. Therefore, it matters what vector is chosen.

The simplest choice is a vector, which vanishes everywhere except on a single point on the lattice. This is a so-called point source. If used, (6.7) gives the result for a propagation from the points x_i to y_i . However, this is as evaluating $\langle \phi(y)\phi(x) \rangle$ for fixed x and y - it gives a propagator from a single point to another single point. It therefore misses the possibility to use translation symmetry to average over the lattice. Therefore, this result will be very noisy. The advantage is that a point source in position space contains in momentum space all modes. Therefore, it is possible to determine the propagator for all momenta by multiplying the vector obtained from conjugate gradient to a vector corresponding to a plane wave for all possible momentum, and thus calculating the propagator in momentum space requires only one inversion, at the expense of being more noisy. On the other hand, inverting Q on a plane-wave vector will give a single point in momentum space, but averages over all lattice points, giving a comparatively good signal-to-noise ratio. However, this requires one inversion per momentum, which is much more expensive. Thus, a trade-off between acceptable noise level, which can also be beaten by the number of configurations, and computing time for the inversions need to be made.

The situation becomes more involved when considering the needs of spectroscopy, as will be detailed more in section 6.8. Calculating a bound-state correlator requires the evaluation of (3.4), build from (3.3). This implies to perform a double sum

$$\sum_{\text{spatial } x,y} \mathcal{O}(x)^\dagger \mathcal{O}(y),$$

which requires to evaluate all the Q from one point to another point, and thus all possible point sources to be multiplied to some point vector are needed, so-called all-to-all

propagators. This requires V inversions and V^2 vector multiplications. This is usually far too expensive. One possibility are so-called wall sources, where not all possible positions are taken into account, but all positions on a hypersurface, assuming that ultimately translation symmetry ensures this to be correct. This reduces the effort to N inversions¹⁴. An alternative is a stochastic evaluation, where the sum is approximated by a number of randomly chosen sources. These issues are beyond the scope of this lecture, but should be considered and researched in the literature before setting out to perform such calculations.

It should be noted that the matrix Q , and thus the fermion propagator, is gauge-dependent. In fact, a single Q^{-1} would be just the fermion propagator. What is not gauge-dependent is the spectrum of Q , and the trace of Q^{-1} , e. g. gives the gauge-invariant, but renormalization-dependent, chiral condensate. This is not a problem if either the configuration is gauge-fixed, or only gauge-invariant combinations of fermion fields are evaluated. In the first case the expressions always have meaning, while in the second the actual gauge does not matter. However, if a configuration is not gauge-fixed, it is in some random gauge¹⁵. This gauge will change during the update, and therefore the fermionic expectation values of two different configurations cannot be compared, if the configurations have not been gauge-fixed before calculating (6.7).

6.8 Hadron spectroscopy

Hadron spectroscopy now amounts almost only to evaluate (6.7) as a result of writing (3.4) in terms of fermion fields. There are just two more issues.

One is that the expression (6.7) is essentially giving only scalars. It is necessary to introduce other quantum numbers.

Consider first the simplest case of some bound state consisting, at the valence level, only out of a fermion and an antifermion, in QCD mesons. To create an operator with the same quantum number is then done as in the continuum by writing an operator

$$\bar{\psi}_i^a(x)\Gamma_{ij}\psi_j^a(x), \quad (6.8)$$

where a is a color index and Γ is a suitable matrix in Dirac and flavor space. E. g. for a π^\pm this would be a tensor product of a $\tau^{1,2}$ in flavor space and a γ_5 in Dirac space, while it would have a τ^3 for the π^0 , but still a γ_5 in Dirac space. This implies that this is just

¹⁴Of course, it could be assumed that the symmetries makes it sufficient to consider just a single pair of points as an estimator, but this will be very, very noisy.

¹⁵Every field configurations defines a gauge, but one which is usually not respecting any symmetries, nor which could be given in terms of any useful gauge condition, lest alone bringing a second gauge field in any straightforward way into the same gauge.

a matrix element of Q^{-1} , which is then contracted with Γ , and ultimately traced over. Thus, technically the same has to be done as in (6.7), just that not only particular matrix elements are needed in position space, but also in color, Dirac, and/or flavor space. In fact, if $\Gamma = 1$, this just returns (6.7).

Likewise, operators for baryons are constructed for a spin 1/2 baryon like

$$\epsilon_{abc}\psi(x)_a^k\Gamma^B(\psi(x)_b^{lT}C\gamma_5\Gamma^A\psi(x)_c^m), \quad (6.9)$$

where $C = i\gamma_2\gamma_0$ is the charge conjugation matrix, $a\dots$ are again color indices and $k\dots$ are flavor indices, and $\Gamma^{A,B}$ are again arbitrary γ matrices. Consider first $\Gamma^A = \Gamma^B = 1$. In this case the term in parentheses forms a scalar, essentially a diquark, which is then coupled to the single remaining spinor to create a spin 1/2 particle. Note that this state is not a parity eigenstate, which is usually an almost conserved quantum number of mass eigenstates, and therefore needs to be projected upon such a state. This is achieved using the projector $(1 \pm \gamma_0)/2$, where $+$ gives the nucleon and $-$ its parity partner, the $N(1535)$. By judicious choices of the $\Gamma^{A,B}$ all other states, which can be created from three quarks, can be constructed.

The antiparticles are created from these operators as usual by hermitian conjugation.

Of course, it is possible to construct much more involved operators, also including relative momenta and more fields, not to mention states like tetraquarks, pentaquarks, and scattering states. They can be used to form a basis for a variational analysis, as discussed in section 4.6.3.

As already discussed in section 4.6.2, such local operators are often not a good approximation, and operators averaged over some lattice sites are better suited. Since expectation values for fermions contain two contributions, the gauge bosons and the fermions, there are two possible ways to perform smearing, just like in section 5.11.

The first one is to smear the gauge fields, using the methods of section 5.3. Then the Dirac propagators Q are evaluated on the smeared links.

The second one is by implementing the smearing itself in the inversion process, and thus in the fermion sector. So far, the inversion was essentially done on some lattice position, or a plane wave. By performing the inversion over a vector which describes a (spatial) extended object, this can be done. After all, the vector on which the inversion is performed has V components. Therefore, the components can be used to create such a structure. Likewise, the final contraction can also be done using such an operator.

The challenge in doing so is that this must maintain gauge covariance. This is not generically the case. One possibility to ensure this is the so-called Jacobi smearing. Given

the point source vector S , the smeared source P is

$$P(x) = \sum_{y_1, y_2, y_3} \left(\sum_{n=0}^N \alpha^n H(x, y)^n \right) S(y)$$

$$H(x, y) = \sum_{\mu=1}^3 (U_\mu(x) \delta_{x+\mu, y} + U^\dagger(x-\mu)_\mu \delta_{x-\mu, y})$$

which is done at fixed time if spectroscopy should be done. The parameters N and α can be tuned independently for each operator to improve, e. g., the signal-to-noise ratio. Physically, the parameter α steers the width of the source, and the final source P can be thought of as something as a proposed wave-function of the hadron. As always, using operators with different parameters can be used to enlarge the operator basis. The links ensure that the final result is gauge-covariant, i. e. transforms like a link between all positions.

Technically, this corresponds to replacing the operator Γ^i in (6.8) and (6.9) by a function of space-time. This also shows why gauge-covariance is important: After all, the hadron should be a gauge-invariant state. This also implies that the plane-wave source of section 6.7 cannot be used for hadronic correlators, as they would not create a gauge-invariant operator.

An alternative to gauge-covariant smearing is to evaluate (6.8) and (6.9), i. e. by fixing the gauge fields using the methods in section 5.5. Then any source can be chosen, as it is then in a fixed gauge. The final averaged result will then be still gauge-invariant as the gauge-dependent parts drop out. Note that this physically corresponds to using operators which are different for different configurations, but only differ in gauge-variant parts.

6.9 Algorithms for fermions

The discussion of the previous sections already point to how to deal with the creation of configurations: Integrate out the fermions, leading to (6.5), and then evaluating the fermion determinant in some numerical way. Since this is then an expression depending on the gauge fields only, it is possible to dice a new gauge field, and then use, e. g., the Metropolis algorithm, to accept or reject the update.

While this procedure works in theory, it is highly inefficient in practice, as this requires to evaluate the fermion determinant for every local update. As the (approximate) calculation of the determinant of a large matrix, even if the matrix is sparse, is very expensive, it appears better to find an update which does not make this necessary.

For this purpose it is useful to introduce a new class of algorithms.

6.9.1 Molecular dynamic algorithms

To introduce this algorithm consider again a single real scalar field. Gauge fields and fermions will be added later.

The basic idea of molecular dynamics is to map the problem of finding a new configuration to a problem of classical mechanics. To this end it consider the Hamiltonian

$$H = \sum_x \frac{p(x)^2}{2} + S, \quad (6.10)$$

where S is the original action of the theory, e. g. (2.23). Thus, to the field at every lattice point an additional momentum $p(x)$ has been introduced. Since by discretization all analytical operations have been turned in sums and differences of the field at various points on the lattice, these momenta are by definition conjugate to the fields at the lattice points, as the action just acts as a potential. Since the lattice has a finite (or at least denumerable infinite) number of points, this is not a field theory, but really a point particle mechanical system. Thus, the system obeys all theorems of classical mechanics, e. g. Liouville's theorem.

Its equation of motion are therefore

$$d_\tau \phi(x) = p(x) \quad (6.11)$$

$$d_\tau p(x) = -\frac{\partial S}{\partial \phi(x)}, \quad (6.12)$$

where again it needs to be emphasized that x is acting as an index, not as a variable. The parameter τ is an additional fictitious coordinate, conceptually just like the Monte-Carlo time, but another, independent variable. Because the classical partition function of this system

$$Z = \int dp(x)d\phi(x)e^{-H}$$

evaluates, after integrating out $p(x)$, to the (discrete) path integral of the original system, the ergodicity hypothesis of classical statistical physics implies that any position $\phi(x)$ is therefore visited with probability $\exp(-S)$, and thus exactly with the required probability.

Therefore, the expectation value of any observable can be computed as a time average in τ ,

$$\langle \mathcal{O} \rangle = \frac{1}{T} \int_{\tau_0}^{\tau_0+T} d\tau \mathcal{O}(\phi(\tau)). \quad (6.13)$$

where τ_0 now takes the role of the thermalization time. Thus, by starting from some random initial field configuration and solving (6.11-6.12) over some time interval T measurements are generated with this method. The necessary initial values for the canonical

momentum can be chosen, in principle, also randomly, but since they only appear in a Gaussian, it can be expected that they will also be Gaussian distributed in equilibrium, and it therefore makes more sense to draw them from the corresponding Gaussian distribution initially.

Of course, there is again no exact analytical solution of (6.11-6.12) in general. Thus, it is necessary to solve the equations (6.11-6.12) numerically using the initial field configuration as initial values, and likewise performing the integral (6.13) analytically. As these are large systems of coupled partial differential equations, this is not entirely trivial.

An analogue to overrelaxation of section 4.3.2 is to regularly replace the actual values of the canonical momenta by new, randomly drawn canonical momenta with again Gaussian chosen values. This induces a drift in a different direction in phase space, without changing that the actual values of the fields are distributed according to Boltzmann weight.

A drawback of this approach is that numerical errors in solving (6.11-6.12) will drive the system out of equilibrium. As this problem increases the coarser the integration steps and the larger the number of equations, this quickly leads to problematic costs in terms of computing time. To avoid this problem, it is possible to combine molecular dynamics with the Metropolis algorithm of section 4.2.

No matter how coarse the integration algorithm, if it is reversible, the result of integrating (6.11-6.12) over a certain interval $\Delta\tau$ can be considered to be a proposed new configuration. Thus, taking this proposal as the input into a Metropolis accept/reject step creates an algorithm which fulfills detailed balance. In this way detailed balance is automatically satisfied, creating a Markov chain. Since it is a global update, it is, however, not possible to just move to the next lattice site. Rather than that a new proposal can then be created by choosing a new, Gaussian distributed set, of canonical momenta, just like in the overrelaxation update.

A convenient choice for solving (6.11-6.12) is the leapfrog algorithm. Time evolution is generated by the Hamilton equation. Using the Hamilton equation, it is found that performing an integration step of $\Delta\tau$ is achieved by first performing

$$p(x) \rightarrow p(x) \tag{6.14}$$

$$\phi(x) \rightarrow \phi(x) + \frac{\Delta\tau}{2}p(x) \tag{6.15}$$

followed by

$$p(x) \rightarrow p(x) - \Delta\tau \frac{\partial S}{\partial \phi(x)} \tag{6.16}$$

$$\phi(x) \rightarrow \phi(x), \tag{6.17}$$

and then once more (6.14-6.15) to finish the update cycle. Because the Hamilton equation is time-reversal invariant, this integration satisfies the preconditions for the Metropolis algorithm. The error is $\mathcal{O}(\Delta\tau^3)$. Performing an evolution of total length τ requires $\tau/\Delta\tau$ steps, and has therefore an error $\mathcal{O}(\tau\Delta\tau^2)$. By tuning τ and $\Delta\tau$ the acceptance probability can be adjusted. Of course, if τ is shorter than the autocorrelation time, which is calculated in τ in the same way as for t_{MC} in section 4.4.2, this will not completely decorrelate the configurations, and the same steps have to be taken as in general.

This can be improved, though not qualitatively but quantitatively, using the Sexton-Weingarten version. This corresponds to multiple updates of types (6.14-6.17), but with shorter time intervals. It starts from (6.16-6.17) with $\Delta\tau/6$, and then alternates with $\Delta\tau/2$, $2\Delta\tau/3$, $\Delta\tau/2$, and $\Delta\tau/6$, for a total of five steps. This can be used to increase $\Delta\tau$ to decrease autocorrelations.

6.9.2 Molecular dynamics for gauge fields

For gauge fields it is convenient to use the links as generalized coordinates. The analogue to (6.10) is

$$H = \frac{1}{2} \sum_{x\mu} \text{tr} \Pi_\mu^2(x) + S,$$

where the canonical conjugated momenta are traceless and hermitian matrices¹⁶. The corresponding molecular-dynamics equations (6.11-6.12) for gauge fields are then

$$\begin{aligned} d_\tau U_\mu &= i\Pi_\mu(x)U_\mu(x) \\ d_\tau \Pi_\mu(x) &= \frac{i\beta}{2d_F} \left(U_\mu(x)V_\mu(x) - (U_\mu(x)V_\mu(x))^\dagger - \frac{1}{d_F} \text{tr} (U_\mu(x)V_\mu(x) - (U_\mu(x)V_\mu(x))^\dagger) \right), \end{aligned}$$

where $V_\mu(x)$ is for the Wilson action just the staples (5.18), where the operation in parentheses takes the traceless, anti-Hermitian part, such that Π is again of the correct type.

These equation can then be used in (6.14-6.17), or in a Sexton-Weingarten-type update. To create random values of the canonical momenta Π , it is possible to decompose $\Pi_\mu(x) = \sum_a \tau^a \pi_\mu^a(x)$, where τ^a are the generators. The values for the real functions π_μ^a can then be selected by Gaussian random values, e. g. of unit width, or any other suitably ergodic choice. Again, this can be adjusted to improve acceptance probability.

Of course, there will be additional terms if matter fields are present. This will now be exemplified for fermions, but is also true for, e. g., the scalars of section 5.11. Note that

¹⁶These are not the conventional canonical momenta for a gauge field, even if the fields instead of the links would be used. After all, the time is a fictitious fifth dimension of a system in equilibrium. Thus, also none of the usual subtleties because of the gauge degrees of freedom arise, nor is gauge-fixing required.

non-trivial topology, i. e. gauge fields like instantons, are in this sense quite problematic. These require large changes, and thus long integration times. But long integration times reduce the acceptance probability for 'normal' gauge field configurations quite substantially. Here again a balance is needed to be both sufficiently ergodic and efficient at the same time.

6.9.3 Updating QCD

Essential all algorithms for fermions in common use today use molecular dynamics, but are still very expensive to do¹⁷. However, this is mainly because of the absence of better working algorithms, not because molecular dynamics is very favorable. After all, since it is a global update, it can make only small changes which lead to long autocorrelation times.

The main reason for this problem arises from the Grassmann fields. Since they cannot be represented numerically, they need to be integrated out. This leads to a determinant, but this therefore couples all points of the lattice. The expression is inherently non-local¹⁸. This prevents the possibility to construct efficient local updates like a heat-bath, as always the complete lattice contributes to the local heatbath. Thus, the determinant would be needed to be recalculated after each local update, making this possibility prohibitively expensive.

As if this problem was not serious enough, also the calculation of the determinant has additional problems. For free fermions, the lowest eigenvalue will be determined by the mass. In QCD, however, there are additional zero modes which stem from the gauge interactions. Depending on the type of fermions, these may exist approximately or exactly on a finite lattice. In addition, there are more small eigenvalues¹⁹, making the conditional number of the Dirac operator unfavorably large for numerical operations. This problem is reduced if the bare quark mass is large.

Finally, the individual eigenvalues are not necessarily positive or even real, depending on the type of Dirac operator. While the determinant itself is necessarily so in the thermodynamic limit, any approximation which does not capture a balanced amount of eigenvalues leads to a violation of the positive weight necessary for the Markov process. If the Dirac operator satisfies at least the so-called γ_5 hermiticity, $\gamma_5 D \gamma_5 = D^\dagger$, this problem

¹⁷The main reason is that the Dirac operator will be needed to be inverted multiple times. However, if observables are build from many quark fields such that many different contractions are necessary, also the calculation of observables can start to become the dominating cost factor. If possible, it is therefore useful to consider whether calculated inverted Dirac operators, the propagators, can be used for multiple purposes in a calculation, or can even be stored if sufficient disk space is available.

¹⁸Which is actually necessary to implement the Pauli principle, in which fermions and bosons differ.

¹⁹This is an effect due to dynamical chiral symmetry breaking in QCD.

is slightly lifted, as then eigenvalues are pairwise complex or real, and thus the determinant is at least real.

If this is the case, two degenerate flavors will yield a guaranteed positive determinant. This is a rather good approximation for QCD. Since the next heavier quark is already quite heavy, the problem for it are much less substantial, and therefore just a degeneracy for the two lightest quarks works comparatively well. It should be noted that there exists Dirac operators, like the so-called overlap operator, which do not have this problem, and are therefore suited even for non-degenerate quarks. However, the numerical calculation of these operators is in itself already very expensive, offsetting this advantage in many circumstances.

Eventually, the resulting algorithms can be optimized in various ways. Each of them has its advantages and disadvantages. Some of them even do not create a perfect Markov chain, but allow for small violations of detailed balance to find a reasonable solution. Here, thus only a single example, suited e. g. for two degenerate flavors of Wilson fermions, will be discussed.

For this purpose so-called pseudofermion fields are introduced to rewrite the determinant again as an exponential

$$\det Q = \int \mathcal{D}\phi^\dagger \mathcal{D}\phi e^{-\phi^\dagger Q^{-1}\phi},$$

which is an non-local action. If there are two degenerate flavors, it follows that $\det Q \det Q = \det Q^\dagger Q$, where γ_5 hermiticity and $\det \gamma_5 = 1$ has been used. The action is then $-\phi^\dagger (Q^\dagger Q)^{-1}\phi$, and therefore it is positive, but still non-local.

Because the theory is still bilinear in the pseudofermion fields, it is actually still possible to determine their distribution exactly. If a field R is created according to a Gaussian of unit width, i. e. determined by $\exp(-R^\dagger R)$, the pseudofermion field $\phi = Q^\dagger R$ will actually satisfy the Boltzmann weight.

With this an update algorithm, the so-called R algorithm, can be constructed.

In a first step the field ϕ is created according to its Boltzmann weight as discussed above. The canonical momenta for the gauge fields are randomly generated. The gauge fields are then updated with this random canonical momentum, while the canonical momenta are updated by

$$\begin{aligned} d_\tau \Pi_\mu(x) &= \frac{i\beta}{d_F} \left(U_\mu(x) V_\mu(x) - (U_\mu(x) V_\mu(x))^\dagger - \frac{1}{d_F} \text{tr} (U_\mu(x) V_\mu(x) - (U_\mu(x) V_\mu(x))^\dagger) \right) \\ V_\mu(x) &= \frac{\partial S}{\partial U_\mu(x)}, \end{aligned}$$

where the action now has contribution from the gauge fields, but also from the term $-\phi^\dagger (Q^\dagger Q)^{-1}\phi$. It thus involves explicitly the pseudofermion fields as well as derivatives of

the (inverse) of the Dirac operator with respect to the links²⁰. The latter can be given as a closed, albeit lengthy, explicit expression involving only the (inverse) Dirac operator and the links. Still, calculating it is the actually expensive step, as the inverse of the Dirac operator is involved. This process is repeated to create full molecular dynamics steps, which are then finally accepted or rejected in a Metropolis update. As with ordinary molecular dynamics, using schemes like leapfrog or Sexton-Weingarten can be used to improve the behavior.

As noted above, this runs into trouble if small eigenvalues occur. One possibility to improve the situation is the so-called Hasenbusch preconditioning. This is done by introducing a second Dirac-operator W , for which already $Q + m_0$, where m_0 is a (larger) mass is a reasonable choice, and a second pseudofermion field with action

$$\det Q^\dagger Q = \int \mathcal{D}\phi^\dagger \mathcal{D}\phi \mathcal{D}\Phi^\dagger \mathcal{D}\Phi e^{-\phi^\dagger (W^\dagger W)^{-1} \phi - \Phi^\dagger ((W^{-1} Q)^\dagger (W^{-1} Q))^{-1} \Phi}.$$

Although now more expensive inversions are needed, the appearing conditions numbers are smaller, and the inversions become (substantially) cheaper.

Note that the pseudofermion fields do no longer play a role after the links have been calculated. All gauge observables involve only the links, and in all fermionic observable, see (6.7), also only the Dirac operators, and thus again only the links U , appear.

Theories which preserve chiral symmetry face one more problem. Because they obey the Atiyah-Singer index theorem, they preserve topology. Thus, the gauge fields configurations separate into sectors with different winding number. There is no continuous evolution connecting different sectors, and therefore molecular dynamics cannot easily change between sectors, yielding an ergodicity problem. This is an important, albeit non-trivial issue, which solution(s) go beyond the scope of this lecture.

²⁰Improved actions, as discussed in section 5.9 for gauge fields, also exist for fermions. Not all of these actions are actually differentiable in the links, especially if particular types of smearing, so-called fat links, are involved. Such cases require more care.

Chapter 7

Finite temperature and density

One of the most important physical phenomena which involves non-perturbative physics is thermodynamics. Because phase transitions imply non-analyticities and separate phases discontinuities it is not possible to treat them in a perturbative approach based on an expansion in a series.

Lattice methods are therefore a suitable alternative. However, separate phases and phase transitions only exist, strictly speaking, in an infinite volume, as was already discussed in the context of quantum phase transitions in section 2.8. This also remains true in thermodynamics. Thus, eventually in a numerical simulation it will always be necessary to perform an extrapolation to infinite volume, again without certainty of whether anything impedes the extrapolation. At the same time the continuum limit is less relevant, as thermodynamic effects are primarily long-range effects. Still, the discretization should be sufficiently good to avoid that discretization artifacts influence the relevant degrees of freedom.

In the following knowledge of basic continuum thermal quantum field theory will be assumed, and it will be concentrated on discussing the changes due to the lattice formulation.

7.1 Finite temperature

As noted in chapter 4, the Euclidean formulation of quantum field theory is actually describing already a quantum field theory in equilibrium at zero temperature and density. In a similar vein as in the reconstruction theorem it is therefore possible to show that any non-equilibrium information is contained in the equilibrium theory, though in practice it may not be possible to extract it with any finite amount of effort.

It is therefore a comparatively trivial step to introduce finite temperature in the system.

For this, recall that a system in equilibrium is static, and a time-dependence no longer exists. Thus, the 'temporal' direction, after all just a conventional choice of one of the equal Euclidean direction, does no longer exist. Importantly, this implies that the spectroscopy of section 3.2 makes no longer sense, as a propagation in time direction no longer takes place. The concept of mass requires refinement.

However, the fourth direction of the lattice is still necessary. According to the Matsubara formalism, it can be used to introduce a temperature into the system. In particular,

$$T = \frac{1}{L_t} = \frac{1}{aN_t}. \quad (7.1)$$

and thus the temperature T is directly related to the temporal extent. This also emphasizes that $L_t \rightarrow \infty$, as was done so far, is equivalent to zero temperature. Note that the Kubo-Martin-Schwinger condition furthermore implies that bosons (including ghosts) need to have periodic boundary conditions, while fermions need to have anti-periodic boundary conditions in the this direction to allow for an interpretation as a temperature¹. There is no condition for the boundary conditions in spatial directions.

This is in principle already all relevant conceptual points when introducing finite temperature in lattice calculations. However, there are many practical issues to be considered.

The first is that also for (7.1) still $a \rightarrow 0$ will be needed to be taken eventually. To keep the temperature fixed this implies that aN_t remains fixed. This is different as so far, since beforehand aN_μ needed to diverge, and thus N_μ had to grow faster than a shrank. This is still true for $L_s = aN_s$, the spatial directions. But not for L_t . This implies that $N_t \neq N_s$, as was the case so far². Thus, the aspect ratio N_t/N_s will be different from one, and necessarily needs to approach zero in the continuum limit. At any finite a this aspect ratio is different from zero. Also, there are different possibilities how to approach zero. This gives rise to an additional systematic error source, which needs to be treated in the same way as in section 4.5 systematic volume and discretization errors. Likewise, this can be considered as a systematic error in treating finite temperature not in an infinite spatial system, where strictly speaking the concept of temperature is not well defined.

The second is that at fixed a temperature can only be changed by changing N_t , and therefore only in fixed step. In particular, there is a maximal temperature, $1/a$. However,

¹Note that this, strictly speaking, implies that this type of boundary conditions also decide whether the case $L_t \rightarrow \infty$ can be interpreted as zero temperature or just as Euclidean field theory. However, since for $L_\mu \rightarrow \infty$ the results become independent of the boundary conditions, this actually plays no role.

²An alternative is the introduction of different discretizations for different directions. If restricted to relatively few lattice sites, this is an advantage, as otherwise the temperature can at fixed a only be changed in steps. However, the implementation of this is then dependent on the theory in question, and will therefore not be discussed here.

as usual, $N_t = \mathcal{O}(1)$ is certainly inviting discretization errors. Likewise, $N_t \approx N_s$ will create large aspect ratio errors. Thus, ideally $1 \ll N_t \ll N_s$ should be ensured, on top of all other conditions discussed in section 4.5. In practice, this is rarely possible.

Another issue is that the introduction of thermodynamics introduces a heatbath. In a relativistic treatment, like quantum field theory, this heatbath has its own four-velocity, u_μ . The standard Matsubara formalism, which has been employed here as well, makes the special choice to work in the rest-frame of the heatbath. This is the reason why the original time-direction is special. There are two consequences of this.

One is that Lorentz symmetry is no longer manifest, as the four velocity of the heatbath is now a distinguished direction. However, Lorentz symmetry is still intact, it is just no longer obvious anymore. By explicitly reintroducing the four-velocity of the heatbath in every step it is possible to regain manifest Lorentz symmetry. While theoretically appealing, this is in practice quite cumbersome and therefore almost never done. While the lattice formulation uses Euclidean space-time to start with, the same considerations apply to the Euclidean $O(4)$ symmetry as well.

The presence of a distinguished four-direction has also another consequence. Any directional quantity, like vectors or tensors, can now be decomposed into elements parallel and perpendicular to this direction. Such quantities will have different properties, which will only merge in the zero temperature limit. In addition, all quantities will depend separately on (pseudo) t and \vec{x}^2 instead of x^4 , making all functions involved. Note that the finite extent in the would-be time direction implies that even in the infinite-volume limit the energy is still quantized in so-called Matsubara frequencies, and zero energy/4-momentum is, because of the Kubo-Martin-Schwinger condition, only possible for bosonic quantities, but not for fermions.

7.2 Finite density and the sign problem

Like temperature, density, or more aptly chemical potential, is another control quantity of a thermodynamic system in equilibrium. Any globally conserved charge can be used to introduce an associated chemical potential to create a system with a non-zero net density of this charge³.

Such a chemical potential can again be introduced using the Matsubara formalism, including the possibility of introducing multiple chemical potentials for different charges.

³Note that electric charge is somewhat subtle in this respect, and any non-Abelian gauge charges immediately create a gauge anomaly when a corresponding chemical potential is introduced, as any such potential explicitly breaks gauge symmetry.

In principle, in the continuum the only thing needed is to add a term

$$\mu N,$$

to the action, where μ is the chemical potential and N is a counting operator for the conserved charge. Naively, this would be just

$$N = \sum_x \phi(x)^\dagger \phi(x),$$

as in the continuum.

However, especially for fermions, it is insufficient to just add such a term.

This is mainly due to the lattice discretization. The number operator is not actually the integrated density, but rather the spatial integral of the fourth component of the Noether current, and thus involves in its definition a direction and thus needs to involve nearest neighbors. A suitable implementation of this feature is then explicitly given by

$$\mu N = \frac{1}{2a} \sum_x (e^{\mu a} \phi^\dagger(x) \phi(x + e_t) - e^{-\mu a} \phi^\dagger(x - e_t) \phi(x)) \quad (7.2)$$

The exponential dependence at leading order generates the correct continuum dependence. Time-reversal invariance then determines that it cannot be just μ , but rather needs to be a function satisfying $f(a\mu) = 1/(f(-a\mu))$ for which the exponential is the simplest, but not the only solution.

Note that the Pauli principle forbids more than one fermion of the same quantum numbers at every lattice site. Thus, for every lattice, there is maximum occupation, and thus maximum possible density.

There is another problem for fermions. The additional term (7.2) is again a bilinear. It therefore becomes part of the Dirac operator. It turns out that for a gauge theory in which the fermion fields are part of a complex representation, which includes QCD, the Dirac operator has no longer a purely real eigenspectrum. This can be seen by the fact that the chemical potential acts like an imaginary gauge-field for such theories. Therefore, the determinant is no longer real (and its square positive). Thus, there is no longer a Boltzmann factor with a probability interpretation, spoiling the usual algorithms to perform simulations.

It should be noted that this is a technical problem, and not a physical one. What actually happens is that all physical quantities are still well-defined. If and how this problem can be solved is yet unclear.

Index

- Action
 - Improved, 98
 - Perfect, 99
- Acton
 - Lüscher-Weisz, 99
- Adaptive algorithm, 54
- APE smearing, 81
- Area law, 97
- Aspect ratio, 64
- Asymptotic freedom, 78
- Autoorrelation time, 57

- Baker-Campbell-Hausdorff formula, 6
- β
 - Abelian, 74
- β -function, 93
- Boltzmann factor, 35
- Bond, 21
- Bootstrap, 56
- Boundary conditions, 8
- Breit-Wigner cross-section, 42
- Brioullin zone, 8, 12
- Bulk phase transition, 26

- Cabibbo-Marinari trick, 84
- Center, 78
- Character expansion, 94
- Chiral theory, 33
- Compact formulation, 74

- Configuration, 35
 - Correlation, 57
 - Drop, 58
- Continuum, 1
- Continuum limit, 23
- Cooling, 81
- Correlation function, 7, 44
 - Connected, 37
 - Time-ordered, 7
- Correlation length, 23
- Correlation time, 58
- Creutz algorithm, 83
- Critical behavior, 26
- Critical exponents, 27
- Critical slowing down
 - Gauge fixing, 88
- Critical slowing-down, 59
- Critical value, 27
- Cross-section, 41
- Cutoff
 - Infrared, 8
 - Ultraviolet, 8

- Density of states, 49
- Derivative
 - Lattice, 10
- Detailed balance, 51
- Disconnected contribution, 37
- Dispersion relation

- Continuum, 16
- Lattice, 16
- Non-relativistic, 16
- Divergence, 19
- Dressing function, 45
- Dynamical exponent, 59
- Energy density, 49
- Energy-momentum relation, 16
- Entropy density, 49
- Equilibrium, 51
- Ergodicity
 - Problems, 62
 - Strong, 51
- Error
 - Statistical, 55
 - Systematic, 60
- Euclidean space-time, 8
- Faddeev-Popov operator, 90
- Feynman gauge, 86
- Feynman rules, 17
 - ϕ^4 theory, 17
- Field configuration, 7
- Field-strength tensor, 74
- Finite-size-scaling analysis, 26
- Form factor, 45
- Fourier transformation
 - Discrete, 44
- Fourier transformation, 12
- Free energy density, 49
- Functional integral
 - Lattice, 11
- Gauge condition, 86
- Gauge field
 - Lattice, 87
- Gauge fixing
 - Los-Alamos method, 87
- Gauge freedom, 74
- Gauge orbit, 75
- Gauge theory, 71
 - Abelian, 72
 - Measure, 74
 - Non-Abelian, 76
 - Non-compact, 71
- Gauge-fixing, 85
 - Numerical, 85
- General relativity, 33
- Ghost field, 89
- Ghost propagator
 - Lattice, 91
- Glueball, 80
 - Charge parity, 81
 - Pseudoscalar, 80
 - Scalar, 80
 - Tensor, 81
- Gluon propagator
 - Lattice, 90
- Graph, 21
- Graph theory, 21
- Gribov-Singer ambiguity, 88
 - Abelian gauge theory, 88
- Ground state, 37
- Haag's theorem, 26
- Heatbath algorithm, 82
- High-temperature expansion, 20
- Hopping expansion, 20
- Hopping parameter
 - Critical, 19
 - Scalar field, 14
- Hypercubic symmetry, 31
- Hypercubic symmetry
 - Representations, 32

- Hysteresis loop, 68
- Improvement, 13
 - Tree-level, 13
- Inertial mass, 16
- Integral
 - Lattice, 11
- Internal energy, 79
- Ising model, 20
- Jackknife, 56
- Jackknife
 - Single-elimination, 56
- Lüscher method, 41
- Λ_L , 93
- Landau gauge, 86
 - Lattice, 87
 - Minimal, 89
- Laplacian
 - Lattice, 10
- Lattice spacing
 - Two-dimensional Yang-Mills theory, 80
- Lattice, 8
 - Extent, 8
 - Hypercubic, 8
 - Hyperrectangular, 8
 - Random, 8
 - Spacing, 8
- Lattice action
 - Free scalar, 11
- Lattice artifact, 60
- Lattice Laplacian
 - Momentum space, 13
- Lattice spacing, 11
- Lattice volume, 8
- Level counting, 40
- Line of constant physics, 25
- Link, 71, 72
 - Adjoint representation, 80
 - Non-Abelian, 77
- Linked-cluster expansion, 22
- Low-energy effective theory, 23
- Markov chain, 52
- Markov process, 52
- Mass
 - Bound state
 - Finite-volume correction, 39
 - Finite-volume correction, 38
 - Ground state, 38
 - Mass zero, 39
- Mass renormalization
 - Scalar, 18
- Master index, 12
- Maxwell theory, 74
 - $\mathcal{O}(a)$ interactions, 76
- Metastability, 69
- Metropolis algorithm, 52
 - Multi-hit, 54
- MOM scheme, 92
- Momentum
 - Continuum, 13
 - Lattice, 8, 13
- Momentum conservation, 13, 17, 46
- Monte-Carlo time, 49
 - Propagation, 53
- Noise, 64
- Operator, 7
 - Fingerprint, 37
- Order parameter, 27, 70
- Overlap, 37, 61
- Overrelaation, 55

- Overrelaxation, 84
- Parallel transporter, 71
- Path integral, 4
- Perimeter law, 97
- Perturbation theory, 15
 - Yang-Mills theory, 90
- Phase shift, 41
- Phase transition
 - Second order, 70
- Phase transition
 - First order, 25, 68
 - Firstorder
 - Signal, 68
 - Second order, 23
- ϕ^4 theory, 14
- Plaquette, 74, 78
- Pole
 - Free particle, 16
 - Residuum, 17
- Polyakov loop, 79
- Probability density, 49
- Projection
 - Zero momentum, 38
- Projection to group, 82
- Propagator
 - Free scalar, 12
 - Momentum space, 13
- Propagator, 44
- Pseudo-experiment, 56
- Quantum fluctuations, 47
- Reconstruction theorem, 9
- Reflection
 - Link, 30
 - Operator, 29
 - Positivity, 29
- Site, 30
- Regularization
 - Lattice, 9
- Renormalization
 - Necessity, 9
- Renormalization trajectory, 25
- Resonance, 40
- Rest mass, 16
- Reweighting, 62
- Rotation symmetry, 31
- Rotational symmetry
 - Violation
 - Estimate, 46
- Roughening transition, 28, 96
- Running coupling, 92
- Scalar QED, 74
 - Lattice, 73
- Scaling limit, 93
- Scattering length, 44
- Scattering state, 40
 - Finite volume correction, 40
- Schwinger function, 29
- Schwinger's line function, 72
- Signal-to-noise ratio, 56
- Simulation, 48
- Smearing, 64
- Specific heat, 26
- Spectroscopy, 63
- Spin, 31
- Spin system, 15
- Staple, 81
- Statistical system
 - Analogy to, 9
- Stochastic process, 50
- String tension, 96
- Strong-coupling limit, 94

- Supersymmetry, 33
- Sweep, 53
- Symmetry, 28, 69
 - Internal, 33
 - Spontaneous breaking, 33
 - Spontaneous breaking, 69
- Systematic error
 - Aspect ratio, 64
 - Discretization, 61
 - Ergodicity, 62
 - Overlap, 61, 64
 - Volume, 61
- Thermalization, 50, 51, 59
- Thermodynamic limit, 9
- Threshold
 - Elastic, 40
 - Inelastic, 40
- Time slice, 6
- Time-slice averaging, 63
- Topological charge, 100
- Topological susceptibility, 100
- Transfer matrix, 6
 - Induced Hamilton operator, 30
- Transition amplitude, 4
- Transition matrix element, 6
- Translation symmetry, 29
- Triviality, 24
- Update, 50
 - Canonical, 54
 - Checkerboard, 53
 - Global, 52
 - Lexicographical, 53
 - Local, 52
 - Microcanonical, 55
- Vacuum expectation value, 37
- Vacuum-to-vacuum transition amplitude,
7
- Variational analysis, 66
 - Preconditioning, 66
- Vertex, 44
 - Bond, 21
 - Four-gluon
 - Lattice, 91
 - Ghost-gluon
 - Lattice, 91
 - Three-gluon
 - Lattice, 91
 - Two-ghost-two-gluon, 92
 - Two-gluon, 90
- Wick rotation, 8
- Width, 42
- Wilson action, 78
- Wilson confinement criterion, 96
- Wilson flow, 65
- Wilson line, 79, 97
- Wilson loop, 79
- Yang-Mills theory, 76
 - Continuum limit, 78
 - Two-dimensional, 80