

Renormalizing Diffusion Models

Jordan Cotler^{1,a} and Semon Rezhikov^{2,3,b}

¹ *Society of Fellows, Harvard University, Cambridge, MA 02138, USA*

² *Department of Mathematics, Princeton University, Princeton, NJ 08544, USA*

³ *Institute for Advanced Study, Princeton, NJ 08540, USA*

^a`jcotler@fas.harvard.edu`, ^b`semonr@princeton.edu`

Abstract

We explain how to use diffusion models to learn inverse renormalization group flows of statistical and quantum field theories. Diffusion models are a class of machine learning models which have been used to generate samples from complex distributions, such as the distribution of natural images. These models achieve sample generation by learning the inverse process to a diffusion process which adds noise to the data until the distribution of the data is pure noise. Nonperturbative renormalization group schemes in physics can naturally be written as diffusion processes in the space of fields. We combine these observations in a concrete framework for building ML-based models for studying field theories, in which the models learn the inverse process to an explicitly-specified renormalization group scheme. We detail how these models define a class of adaptive bridge (or parallel tempering) samplers for lattice field theory. Because renormalization group schemes have a physical meaning, we provide explicit prescriptions for how to compare results derived from models associated to several different renormalization group schemes of interest. We also explain how to use diffusion models in a variational method to find ground states of quantum systems. We apply some of our methods to numerically find RG flows of interacting statistical field theories. From the perspective of machine learning, our work provides an interpretation of multiscale diffusion models, and gives physically-inspired suggestions for diffusion models which should have novel properties.

Contents

1	Introduction	3
1.1	Preliminaries	3
1.2	A sketch of bridging the physics and ML perspectives	6
1.3	History	10
1.4	Summary of this paper	10
2	Variational inference and sampling	12
2.1	Basic objects	12
2.2	Variational lower bounds	13
2.3	Sampling algorithms	14
3	Review of mathematical aspects of Diffusion Models	17
3.1	Score-based generative modeling	17
3.2	Taking the continuum limit	21
4	Lattice field theory and the renormalization group	24
4.1	Lattice discretization	24
4.2	Renormalization group	25
4.3	A toy analog of the renormalization group	28
4.4	Interfacing the renormalization group with modeling and simulation	28
4.5	Comments on RG fixed points and phase transitions	29
4.6	Continuous-time formalism	31
4.7	The concept of Effective Field Theory	34
5	Renormalizing diffusion models and multiscale modeling	37
5.1	Overview	37
5.2	Variationally optimizing the proposal distribution: normalizing flows	39
5.3	Learning the renormalization group	40
5.4	Variational inference for sampling and diffusion models	43
5.5	More history	46
5.6	Other RG schemes	48
6	Finding ground states of quantum field theories	56
6.1	Review of difficulties with variational methods in quantum field theory	56
6.2	Real-valuedness of the ground state wavefunction	58

6.3	Learning ground states of QFTs with diffusion models	58
7	Numerically learning RG flows	61
7.1	Overview	61
7.2	Conventions and estimators	62
7.3	The Carosso and Polchinski flows	63
7.4	Numerics	64
8	Discussion	68
A	A brief review of literature on diffusion models	69
B	Lattice discretization of functional derivatives	70
C	Review of the Exact Renormalization Group	71

1 Introduction

1.1 Preliminaries

Generative machine learning algorithms try to learn a complex, high-dimensional distribution

$$\mu = p(x) dx, \quad x \in \mathbb{R}^D, \quad (1.1)$$

by finding optimal parameters θ^* for a family of distributions

$$p_\theta(x) dx, \quad (1.2)$$

which is defined either explicitly in terms of a density $p_\theta(x)$ or implicitly in terms of a sampling scheme. Typically the distribution $p(x)$ to be learned is in a sense ‘naturally occurring’, for instance a distribution over natural images. For many such distributions, each sample x organizes into a *field*; for example, an image is a function $\vec{\phi} : \{1, \dots, N\} \times \{1, \dots, N\} \rightarrow \mathbb{R}^3$, and $x = \{\vec{\phi}_{i,j}\}_{i,j=1}^N$ are the values of ϕ on a lattice embedded in \mathbb{R}^2 . Here each lattice site can be viewed as a pixel. Thus, in the case of images, a distribution μ_N over such $\vec{\phi}$ ’s may be viewed as a finite-dimensional approximation to a hypothetical infinite-dimensional distribution

$$\mu = \frac{1}{Z} P[\vec{\phi}] \mathcal{D}\phi = \frac{1}{Z} e^{-S[\vec{\phi}]} \mathcal{D}\vec{\phi}, \quad \vec{\phi} : \mathbb{R}^2 \longrightarrow \mathbb{R}^3, \quad (1.3)$$

where a fixed $\vec{\phi}(y) = (\phi_1(y), \phi_2(y), \phi_3(y))$ for $y \in \mathbb{R}^2$ provides an image at infinite resolution. For simplicity, it is easier to think about black and white images which are naturally described by a single field $\phi(y)$ and would have an associated hypothetical infinite-dimensional distribution

$$\mu = \frac{1}{Z} P[\phi] \mathcal{D}\phi = \frac{1}{Z} e^{-S[\phi]} \mathcal{D}\phi, \quad \phi : \mathbb{R}^2 \longrightarrow \mathbb{R}. \quad (1.4)$$

In statistical field theory, one studies infinite-dimensional distributions akin to the one above, where the field $\phi(y)$ for $y \in \mathbb{R}^2$ corresponds to e.g. the magnetization of a magnet at position y , the temperature of a metal at position y , the density of a plasma at position y , etc. Then the ‘action’ $S[\phi]$ would be derived from a physical model. For example, a model called “scalar ϕ^4 theory” which has action

$$S[\phi] = \int_{\mathbb{R}^2} dy \left(\frac{1}{2} \nabla\phi \cdot \nabla\phi + \frac{1}{2} m^2 \phi^2 + \frac{\lambda}{4!} \phi^4 \right) \quad (1.5)$$

is central to the study of statistical properties of magnets.

In lattice field theory, one constructs computational methods for sampling from finite-dimensional analogs of distributions like $P[\phi] = \frac{1}{Z} e^{-S[\phi]}$ for $S[\phi]$ in (1.5) in order to compute expectation values of observables of interest, such as the average magnetization of a magnet over certain lattice sites. Sampling from such distributions is challenging due to their multi-modal and high-dimensional nature, and there have been attempts to apply machine-learning techniques, e.g. flow-based generative modeling, to improve sampling [1–4]. In the related setting of quantum field theory, one wishes to gain access to expectation values in the quantum *ground state*, which in certain physically relevant cases turns out to be representable by a non-negative distribution (see Section 6.2). In this quantum ground state setting an explicit form of the density is generally unavailable, but one can search for the quantum ground state via a minimization problem directly analogous to the problem of variational inference [5].

Due to the high-dimensional nature of all of the above problems, one is required make approximations, and to incorporate prior knowledge about the distribution $p(x)$ or $P[\phi]$ into the structure of the model. In generative machine learning, it has been found that one can successfully learn high-dimensional distributions by leveraging a stochastic differential equation (SDE)

$$dx = f(x) dt + \sigma(x) dB_t \quad (1.6)$$

which induces an on evolution $p_t(x)$ from the desired complex distribution $p_0 = p$ to a simpler distribution p_∞ , which is typically Gaussian or uniform. There is an associated *inverse* SDE

$$dx = (-f(x) - \sigma^T \sigma(x) s(x)) dt + \sigma(x) dB_t \quad (1.7)$$

where $s(x)$ is the *score function*

$$s(x) = \nabla p_t(x), \quad (1.8)$$

which induces a flow from p_∞ to p_0 . In *diffusion modeling*, one variationally parameterizes the *score function* $s(x)$ via a neural network with parameters θ , so that $s_\theta(x)$ is used as a proxy. For a fixed θ , one uses the approximate score function and the inverse SDE to generate approximate samples from p_0 given easy samples from p_∞ . By minimizing a functional computed in terms of these samples, one can improve upon the initial guess of θ so as to make $s_\theta(x)$ closer to the true $s(x)$, and thus improve the sampler. Here, implicit priors about the distribution $p(x)$ are encoded

into the diffusion SDE, the construction of the neural network s_θ , and the particular details of the training scheme, which involves e.g. a choice of numerical SDE solver.

In the statistical physics of fields, it is well-established that the *renormalization group* provides important analytical and numerical insights into the study of the pertinent distributions [6]. The basic conceptual insight is that our description of a physical system is contingent on the precision of our measurement apparatus. In particular, if we can only probe a system in a coarse-grained manner which is insensitive to system properties that are sufficiently small, then our *effective* description of the system can neglect those short-distance properties. A useful conceptual example is a fluid; if your measurement apparatus cannot resolve individual atoms, then the equations you use to describe what you can measure about the fluid need not model the individual atoms but rather may treat the system as a continuum. At a more practical level, the question the renormalization group seeks to answer is: given a probability distribution $P[\phi] = \frac{1}{Z} e^{-S[\phi]}$ over fields $\phi(y)$, how do we find a *coarse-grained* description of the distribution adequate to accurately reproduce expectation values of observables smoothed over large distances? We emphasize that for various physical models of interest, fundamental notions such as phase transitions and critical exponents are studied or even defined in terms of the renormalization group, and a wealth of physical intuition is available about the behavior of the renormalization group in the context of specific models (see [7] for a modern treatment of the subject).

Since the renormalization group instantiates a type of coarse-graining process, it may come as no surprise that it can be described in terms of a heat-flow-like equation for probability densities on the space of fields, and also via a stochastic-differential equation for the field variables ϕ . There are different choices of *renormalization group scheme*, corresponding to different ways of implementing the coarse-graining process; essentially all reasonable continuous-space renormalization group schemes are characterized by their associated SDE, as we show in Section 4. To illustrate, consider the *Carosso scheme* [8]

$$d\phi_t(y) = \Delta\phi_t(y) dt + dB_t(y), \quad (1.9)$$

where t is a ‘time’ that parameterizes how much we have coarse-grained our system, Δ is the (spatial) Laplacian, and $dB_t(y)$ describes a t -dependent Gaussian random field. In fact (1.9) induces a deterministic flow $P_t[\phi]$ where $P_0[\phi] = P[\phi]$. We will elaborate on the details and meaning of these equations later on, but for now it suffices to say that $P_t[\phi]$ for $t > 0$ provides our desired coarse-grained description of the physical system described by $P[\phi]$.

How should we conceive of the coarse-grained probability functional $P_t[\phi] = \frac{1}{Z_t} e^{-S_t[\phi]}$? Going back to the action in (1.5), a t -dependent version would look like

$$S_t[\phi] = \int_{\mathbb{R}^2} dy \left(\frac{a(t)^2}{2} \nabla\phi \cdot \nabla\phi + \frac{1}{2} m(t)^2 \phi^2 + \frac{\lambda(t)}{4!} \phi^4 + \dots \right) \quad (1.10)$$

where the mass m and λ are now functions of t , along with a new multiplicative factor a which sets the scale of the kinetic term, and the \dots terms represent new kinds of interaction terms that we pick up in the process of the flow. The functions $a(t)$, $m(t)$, and $\lambda(t)$ (as well as those suppressed

in the \dots terms) are quantities of physical interest, since they capture information about what we can actually measure with a particular precision of our measurement apparatus. They comprise the t -dependence of the Taylor coefficients of the log-density of $\mu_t = P_t[\phi] \mathcal{D}\phi$. Terminologically, the distribution $\mu_0 = P_0[\phi] \mathcal{D}\phi$ is called the *UV distribution*, while the distribution $\mu_\infty = P_\infty[\phi] \mathcal{D}\phi$ is called the *IR distribution*. These names derive from the fact that UV light is short-wavelength, and IR light is long-wavelength; accordingly the UV distribution captures short-distance physics and the IR distribution captures long-distance physics.

Putting some of the above ideas together, if we wish to sample from the distribution associated to a statistical field theory, we can use a finite-dimensional discretization of the SDE associated to a renormalization group scheme to define a diffusion model which *learns the inverse SDE to the renormalization group flow*. In this paper we will describe in detail how to design such models. Crucially, many of the parameters of these models have a direct physical interpretation, a property that is difficult to achieve in typical black-box applications of machine learning methodologies to physics. In particular, we will show how one can *compare* results from different choices of models in this class which are based on different RG schemes (i.e. different SDEs) by deriving appropriate ways of renormalizing the fields (i.e. rescaling model variables) such that comparisons can be made. We believe that the physical interpretability of these machine learning models will make debugging and modifying such models more straightforward in physical applications, and possibly inspire new techniques in more general machine learning domains. We present empirical results for this class of models, and more broadly aim to make clear the precise connection between diffusion models and the renormalization group for the benefit of physicists and machine learning practitioners, in particular so that ideas and techniques can be communicated between these different communities of researchers.

1.2 A sketch of bridging the physics and ML perspectives

From the perspective of lattice field theory, a physicist might have invented diffusion modeling as follows. The distributions μ_0 pertinent for, say, studying quantum chromodynamics, are complicated and often forced to be multimodal due to the existence of physical quantities such as the *topological charge* of a gauge field, which contribute to *critical slowing down* of direct MCMC sampling of μ_0 [9]. A well-known MCMC technique is *bridge sampling* or *parallel tempering* [10–12]: one runs parallel MCMC samplers for a family of distributions μ_t , $t = 0, \dots, T$ where the *bridge* μ_T is a simpler (“high-temperature”) distribution which may be unimodal or approximately Gaussian. See Figure 1 for an illustration. The method works by using samples from each μ_t as proposals for each μ_{t-1} , with acceptance ratios chosen such that the joint Markov chain stills samples from $\mu_0 \times \dots \times \mu_T$. Because μ_T is simple, e.g. unimodal, the sampler for μ_T mixes quickly, which in some cases improves the overall mixing of the sampler from μ_0 , despite the additional overhead cost of running several parallel samplers.

A key problem is engineering a natural sequence of μ_t ’s which form a good bridge. Fortunately, physics gives us a *preferred bridge*: the flow of μ_0 given by the renormalization group (see e.g. [13,

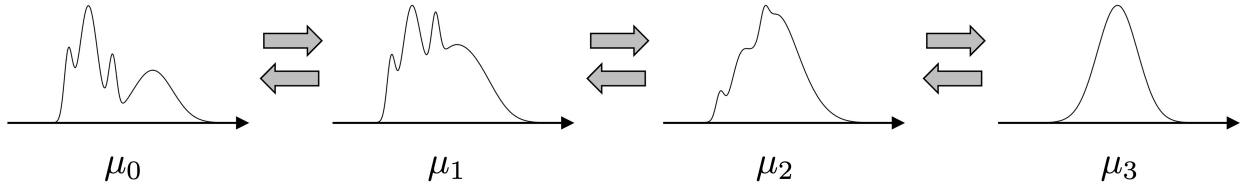


Figure 1: *Parallel tempering*. It is difficult to sample from multimodal distributions like μ_0 , due to slow mixing times of Markov chains. A common strategy is to connect the distribution of interest to a simple unimodal distribution from which sampling is straightforward, and then use coupled MCMC chains (or diffusion models, or normalizing flows) to speed up sampling from μ_0 . The process connecting μ_0 to a simple distribution μ_3 is often a noising or averaging process, and can sometimes be given a physical meaning in the setting of the renormalization group. Multimodality of μ_0 is in physical settings often associated with phase transitions or approximately conserved quantities like the topological charge.

14]). Unfortunately, the probability densities of the μ_t (the images of μ_0 under RG flow) are not known explicitly; this presents a problem since to construct an MCMC sampler for a μ_t , one needs to know the density of that μ_t somewhat explicitly. One can make approximate guesses for the densities of μ_t using analytical techniques, and indeed such methods can be used to improve sampling (see e.g. [15]). Alternatively, one can simply *parameterize* the densities of the μ_t (or at least their score functions, which are sufficient to construct MCMC samplers) and instead perform a *variational, adaptive variant of bridge sampling* where one *learns* the samplers for the distributions in the bridge. Implementing this physically natural idea would lead directly to the models described in this paper.

However, the diffusion SDEs associated to renormalization group schemes have correlations between model variables, in contrast to the original SDEs used in the machine-learning community for image modeling. Indeed, discretizing the Carosso scheme (1.9) on a 2D lattice, we get the SDE

$$\begin{aligned} d\phi_{i,j} &= (\phi_{i-1,j} + 2\phi_{i,j} - \phi_{i+1,j}) dt + (\phi_{i,j-1} + 2\phi_{i,j} - \phi_{i,j+1}) dt + (dB_t)_{i,j} \\ &= (\Delta\phi)_{i,j} dt + (dB_t)_{i,j} \end{aligned} \quad (1.11)$$

where here Δ is the discrete Laplacian. The above can be written more schematically as

$$d\phi_{i,j} = f(\phi_{i,j}, \phi_{i-1,j}, \phi_{i+1,j}, \phi_{i,j-1}, \phi_{i,j+1}, t) + (dB_t)_{i,j} \quad (1.12)$$

for the variables $\{\phi_{i,j}\}_{i,j}$. In contrast, diffusion models commonly use SDEs where all variables diffuse independently [16, 17], namely

$$d\phi_{i,j} = f(\phi_{i,j}, t) dt + (dB_t)_{i,j} \quad (1.13)$$

for some variables $\{\phi_{i,j}\}_{i,j}$. The key difference between (1.12) and (1.13) is that in the former equation the dynamics of the variables are coupled via the SDE, whereas in the latter equation each variable evolves according to its own decoupled SDE.

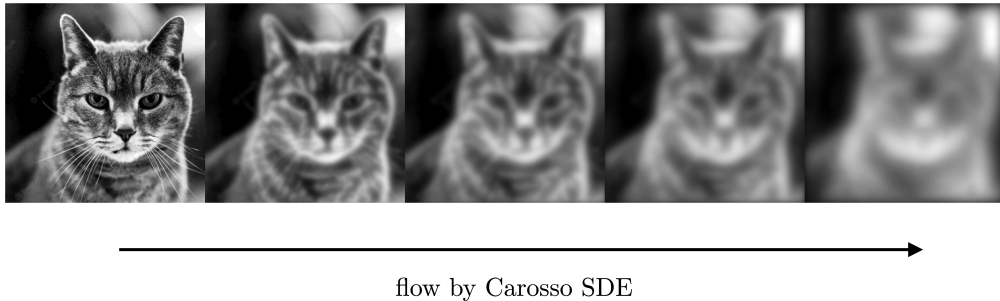


Figure 2: Evolving an image of a cat using the Carosso SDE.

In the discretized Carosso SDE (1.11), the variables diffuse such that high-frequency fluctuations of the discretized field $\phi_{i,j}$ (e.g. pixels of an image) will be turned to data-independent noise. As the evolution continues, progressively lower frequencies will in turn be transformed into data-independent noise; see Figure 2. Thus, in the inverse process, low-frequency modes of the field are generated before the high-frequency modes, which may seem intuitively desirable. Remarkably, the image-modeling community has already studied precisely the SDE (1.11) for this reason [18], and in general, many *multiscale* diffusion models have been developed [19–22] and found to improve computational efficiency and state-of-the-art performance. Even in commercial image-generation systems such as Midjourney, one generates high-resolution image models by training a hierarchy of image-conditioned diffusion models which first generate low-resolution images and then repeatedly upsample the images to higher resolution.

We suspect that various architectures of multiscale diffusion models correspond rigorously to implementations of various renormalization group schemes, and in turn the physical understanding of renormalization group schemes may give rise to new designs and diagnostic techniques in machine learning. From the perspective espoused here, the active use of multiscale diffusion modeling in machine learning suggests that the generative modeling community has independently understood various computational advantages of the renormalization group, and that there are more insights to be mined in this direction.

We also argue that ideas from diffusion modeling in machine learning can be imported into numerical methods for physics problems. For a typical physical system, one can write down an ansatz for a *microscopic action* or *Hamiltonian*, as well as natural *coarse-graining* procedures (also called *renormalization group flows*). These procedures define *implicit* intermediate-scale models, about which we may know certain pertinent parameters but for which we typically do not know an explicit action or Hamiltonian. In a number of cases, coarse-graining procedures can be chosen to make the simulation of an estimated coarse-grained system *simpler*. Thus, at a very high level, one has two processes that one can run in parallel:

1. Simulate the microscopic system, and use the generated dynamics to improve estimates for the dynamics of coarse-grained versions of the system (see Figure 3(a)); and

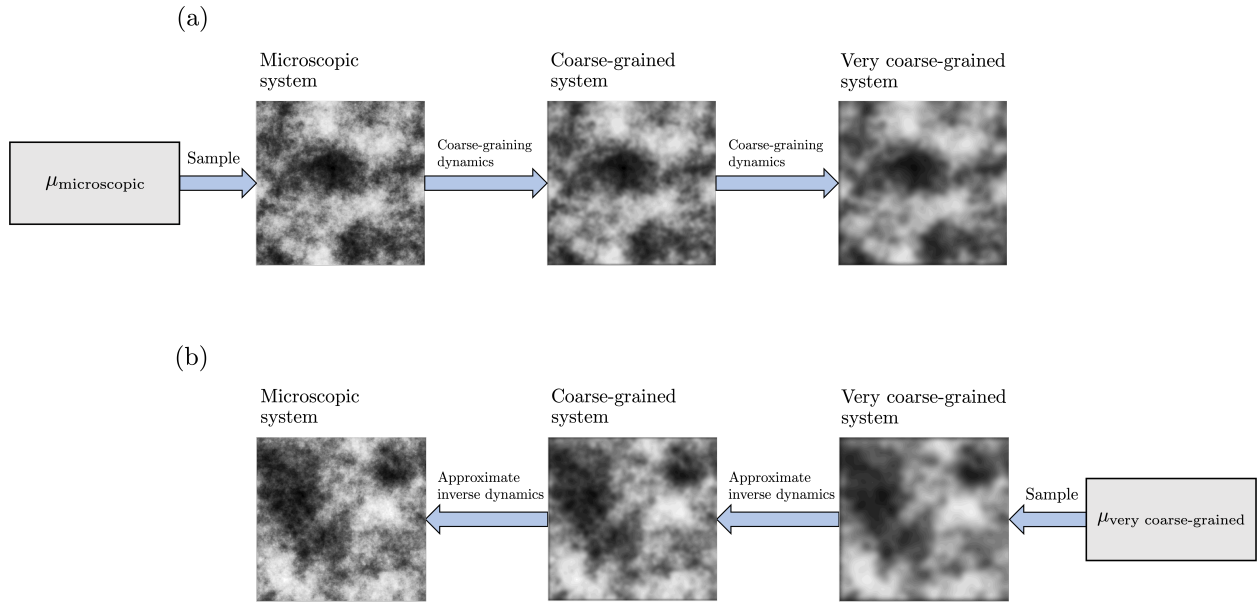


Figure 3: Physical systems tend to be describable by models of varying levels of accuracy (often associated with a hierarchy of length scales). Fine-grained models are accurate but rarely tractable, while coarse-grained models are efficient but have unknown error terms. (a) In the renormalization group framework, it is possible to instantiate coarse-graining dynamics on a microscopic model to obtain quantifiable coarse-grained systems. (b) It is also possible to use coarse-grained systems in conjunction with an approximate inverse to coarse-graining dynamics to improve the efficiency of simulations of a microscopic system. Developing a machine learning toolkit to improve coarse-grained physics models in a way that is physically interpretable and justified is a significant problem at the intersection of physics and computer science. This paper studies the case of simple statistical field theories, using the connection between the renormalization group and diffusion models, but the problem is not limited to this setting.

2. Simulate coarse-grained versions of the microscopic system (more efficiently than the corresponding microscopic system), and use the resulting samples to improve the determination or simulation of the microscopic system. This is done by learning an *approximate inverse* to the coarse-graining procedure (see Figure 3(b)).

This schematic methodology could theoretically be applied to many modeling problems, such as those in condensed matter, high-energy physics, quantum chemistry, and plasma physics. We hope that the theoretically clean connection between diffusion modeling and the renormalization group for scalar field theories presented in this paper can be expanded (by taking non-relativistic limits, incorporating fermions and gauge symmetries, and elaborating on the connections to Effective Field Theory) to create a robust, physically-interpretable machine learning toolkit for numerical physics problems. We will begin to develop such a toolkit in this paper.

1.3 History

The present paper posits a systematic identification between multiscale diffusion models and the renormalization group. The ideas behind diffusion modeling in machine learning arise from several different, interrelated perspectives [23–25], in part motivated by the idea of the renormalization group [23]. These approaches were later unified [16, 17] and organized according to a class of SDEs where the model variables diffused in an *uncorrelated* fashion, as in (1.13), and thus do not implement renormalization group transformations in the sense explained in this paper. Later, variants of multiscale diffusion modeling were proposed by many authors [18, 20–22, 26] and the topic remains an active area of research.

On the other hand, in lattice field theory, several lines of work were inspired by Lüscher’s pioneering paper [27], which pointed out that when trying to sample from lattice $SU(N)$ gauge theory, which is a distribution μ over a compact manifold, there is an implicit characterization of a flow f_t called the *Wilson Flow* such that $(f_t)_*\mu$ limits to a uniform distribution (i.e. a Haar measure). In fact, this flow ends up being a lattice analog of the gradient flow of the Yang-Mills functional [28]. As such, it is a PDE for the fields with a gradient flow interpretation which has as its highest-order term the Laplacian operator acting on the fields. It was noticed that this flow smooths out the fields and gives rise to an operation akin to the renormalization group. This led to a significant amount of work applying the ‘gradient flow’ to lattice field theory methods which has had many applications, e.g. to set the scale of lattice QCD simulations [29]. The precise meaning of the ‘gradient flow’ was eventually clarified by Carosso [30], who explained that there is a valid renormalization group scheme which is a stochastic analog of the gradient flow. In this paper we will take the Carosso scheme [8] as a starting point and use it to illustrate the connection between diffusion modeling and the renormalization group since the scheme has a simple lattice implementation (which indeed turns out to have been studied in the generative modeling community [18]). We will then turn to more general schemes.

In a different line of work following Lüscher’s idea, there are a series of papers which apply machine learning methods to attempt to *learn* a trivializing flow in order to sample more efficiently from UV distributions [2, 31, 32]. These works can be viewed as part of a broader trend to apply black-box machine learning methods to problems in field theory, molecular dynamics, and condensed matter [33]. Flow-based generative modeling methods have demonstrated improved integrated autocorrelation statistics for samplers [2] and improved exploration of topological charge sectors [31], but evaluation of their effectiveness at larger lattice sizes is still ongoing [32]. Crucially, trivializing flows are *not unique*, and in these ML-based methods, the flow learned *has no physical interpretation*, making it more challenging to interpret, debug, and reason about the models.

1.4 Summary of this paper

In this paper we propose a framework for multiscale diffusion and flow-based generative modeling techniques, and then specialize to problems in statistical field theory in which the models *learn the*

inverse renormalization group flow of a field theory. We begin in Section 2 with a basic review of Bayesian inference and Monte Carlo sampling. In Section 3, we review the mathematical ideas and basic implementation details of diffusion models. In Section 4, we review the basic ideas of lattice field theory, and outline how the renormalization group is required to perform calculations of physical quantities of interest. In Section 5, we explain how to design a diffusion model based on the Carosso scheme and how to use it to sample from the UV distribution of a lattice scalar field theory. Importantly, we will explain how the renormalization group requires that physical quantities are extracted from *rescalings* of the variables in this model. This idea may be of particular interest to researchers in machine learning. In *multiscale* diffusion models, the corruption process affects different frequencies at different rates, and the ideology of the renormalization group suggests to *appropriately* rescale all variables to *focus in* on the “interesting” part of the distribution. In particular, there is an associated diffusion process for the *renormalized* fields which has qualitatively different behavior from typical diffusion processes, because it can have fixed points which are highly non-Gaussian (although a generic distribution flows to a Gaussian fixed point). For multiscale diffusion models associated to renormalization group schemes, one can derive an approximation to the SDE for the renormalized fields, which can be used in physical modeling applications and may be helpful in applications to generative modeling.

It is a fundamental fact of statistical field theory that certain quantities of physical interest, such as critical exponents, are independent of the choice of renormalization group scheme. Later on in Section 5, we write down the SDEs for general Wegner-Morris exact renormalization group schemes. As a special case, we describe a diffusion model based on the more traditional Polchinski exact renormalization group scheme [34]. Next we explain how renormalization group ideas indicate how to rescale variables to meaningfully compare results between the schemes. We provide formulae and intuitions for how different multiscale diffusion SDEs should be compared on an equal footing, which are results of potential interest to the ML community.

In Section 6 we explain how to use diffusion modeling to provide an explicit class of variational methods for learning the ground states of quantum field theories. The basic challenge for variational methods for computing ground states is that, as noted by Feynman [35], without an inspired parameterization, the gradients of the variational problem blow up because they are primarily sensitive to unimportant high-frequency fluctuations. In low-dimensional quantum field theories, variational methods based on the renormalization group (e.g. DMRG [36], MPS [37], and MERA [38]) have been helpful for computing quantities of physical interest like correlation functions and critical exponents (see e.g. [39]). However, these methods have not yet been successfully adapted to general higher-dimensional field theories, in part due to computational limitations. The multiscale-diffusion-based ansatz presented in this paper is a novel, physically grounded variational method for finding the renormalization group flow of the ground state of a quantum field theory, and we provide some intuition about why the exploding gradients problem observed by Feynman may be less significant for these methods.

In Section 7, we show tests of a numerical implementation of one of our renormalization group-

based algorithms for sampling from statistical field theories. Our results show that flow-based methods optimized with objectives derived via the renormalization group can learn the RG flows of basic scalar field theories such as ϕ^4 theory.

We end the paper in Section 8 with a discussion of directions for further research, including subjects that might benefit from being studied using multiscale diffusion models.

In Appendix A, we provide a brief review of literature on diffusion models. In Appendix B, we describe some technical points about the lattice discretization of functional derivatives that we utilize in the body of the text. Finally, in Appendix C, we give a short review of the Exact Renormalization Group.

2 Variational inference and sampling

2.1 Basic objects

As we recall from the Section 1, in typical modeling problems we are confronted with a distribution of interest

$$p(x) dx = \frac{1}{Z} e^{f(x)} dx, \quad x \in \mathbb{R}^D, \quad (2.1)$$

where Z is a normalization constant such that $\int dx p(x) = 1$. We may have access to this distribution for a finite number of empirical samples $\{x_i\}_i$, or through an expression for the log-density $f(x)$; typically, the normalization constant Z is unknown and challenging to compute.

In variational inference, we try to approximate the distribution of interest p by a variational family of distributions

$$p_\theta(x) dx = \frac{1}{Z_\theta} e^{f_\theta(x)} dx, \quad x \in \mathbb{R}^D, \quad (2.2)$$

which we may represent explicitly via $f_\theta(x)$ or implicitly via a method to sample from $p_\theta(x)$. A convenient quantity one can use to represent $p_\theta(x)$ is the *score function*

$$s_\theta(x) = \nabla_x \log p_\theta(x) = \nabla_x f_\theta(x). \quad (2.3)$$

The score function is a vector-valued function that points in the direction of increasing probability mass. The score function has many convenient properties, including that it is independent of the normalization constant Z_θ , and that many sampling algorithms only require knowledge of the score function to draw samples from $p_\theta(x)$. Models that represent $p_\theta(x)$ via its score function are called *score-based* models.

To find the optimal parameter θ in the variational family, we must optimize an objective function depending on θ . Many such functions with convenient properties can be built out of the *KL divergence*

$$\text{KL}(p|q) = \int dx p(x) \log \left(\frac{p(x)}{q(x)} \right) = -H(p) - \mathbb{E}_{x \sim p(x)} [\log q(x)]. \quad (2.4)$$

Here $H(p) := -\int dx p \log p$ is the *entropy* of p . The KL divergence is non-negative and only zero when $p = q$. A related quantity is the so-called Fisher divergence¹

$$D_F(p|q) = \int dx p(x) |\nabla \log p(x) - \nabla \log q(x)|^2, \quad (2.5)$$

which is the basis of the *score matching* techniques reviewed in the next section.

Most natural objective functions that can be used for variational inference involve an intractable integral over x , and thus require *gradient estimators* for an optimal θ to be found using (stochastic) gradient descent. For example, if we choose the objective function to be $L_1(\theta) = \text{KL}(p_\theta|p)$, then

$$\nabla_\theta \text{KL}(p_\theta|p) = \mathbb{E}_{x \sim p_\theta} [\nabla_\theta f(x)(f_\theta(x) - f(x) + 1)] \quad (2.6)$$

where we have used the identity $\nabla p = p \nabla \log p$. On the other hand, if we choose the objective function to be $L_2(\theta) = \text{KL}(p|p_\theta)$, then

$$\nabla_\theta \text{KL}(p|p_\theta) = -\mathbb{E}_{x \sim p(x)} [\nabla_\theta f_\theta]. \quad (2.7)$$

Thus to optimize $L_1(\theta)$ using stochastic gradient descent we need access to samples from p_θ , whereas to optimize $L_2(\theta)$ we need access to samples from p . Moreover, to optimize $L_1(\theta)$ we need access to $f(x)$ and $f_\theta(x)$, which may or may not be possible depending on the model architecture and the nature of the problem; to optimize $L_2(\theta)$ we need access to $\nabla_\theta f_\theta$ which may also be challenging for certain score-based models. The usage of the backwards or forwards KL divergences $\text{KL}(p_\theta|p)$ and $\text{KL}(p|p_\theta)$ in variational inference tends to have different tradeoffs due to their zero-forcing and zero-avoiding behavior respectively [43], although each version can be useful in practice. Moreover, there are a wide variety of methods for writing gradient estimators in certain classes of models which may improve the variance of the estimators [44].

2.2 Variational lower bounds

One can use variational inequalities to produce new optimization objectives involving new variables. Introducing a new random variable z with probability density $r(z)$, we have by Bayes' rule

$$p_\theta(x) = \int dz p_\theta(x|z) r(z). \quad (2.8)$$

In many cases, z can be chosen in a meaningful way, such that the form of $p_\theta(x|z)$ and $r(z)$ can be taken as providing a *definition* of the variational family $p_\theta(x)$. Then for any auxiliary distribution $q(z|x)$, by multiplying and dividing by q one has

$$p_\theta(x) = \mathbb{E}_{z \sim q(z|x)} \left[\frac{p_\theta(x|z)}{q(z|x)} \right] \quad (2.9)$$

¹There are several related but not identical quantities involving expectation values of squares of gradients, including the Fisher information $\int dx p(x|\theta)(\partial_\theta p(x|\theta))^2$. The Fisher divergence of (2.5) is connected to the KL divergence by de Bruijn's identity (see [40, Definition 2.5]) and by the Cramér-Rao bound in the case of location parameters θ (see [41, Section 17.7]), and possesses entropy-like monotonicity properties under taking i.i.d. sums [42].

and thus by Jensen’s inequality we have the *variational bound*

$$\log p_\theta(x) \geq \mathbb{E}_{z \sim q(z|x)} \left[\log \left(\frac{p_\theta(x|z)}{q(z|x)} \right) \right]. \quad (2.10)$$

In using Jensen’s inequality we have thrown away the term $\text{KL}(q(z|x)|p_\theta(z|x))$, which can thus be thought of as an error term.

One can think of (2.9) as the justification for *importance-sampling*: given knowledge of $p_\theta(x|z)$, one can estimate (or *define*) $p_\theta(x)$ by sampling x with some importance weighting distribution q . Now, by (2.4), one sees that optimizing $\text{KL}(p|p_\theta)$ is equivalent to optimizing $\mathbb{E}_{x \sim p}[\log p_\theta(x)]$. Unfortunately, Jensen’s inequality implies that the importance-sampling estimator for $\log p_\theta(x)$ is biased downwards, with an error term proportional to how well $q(z|x)$ approximates $p_\theta(z|x)$. In variational inference one typically optimizes the variational lower bound of (2.10) (i.e. the right-hand side) instead of $\mathbb{E}_{x \sim p}[\log p_\theta(x)]$, using judicious choices of $q(z|x)$ which may themselves be chosen variationally from a family q_θ .

2.3 Sampling algorithms

After producing an estimate for the log-probability density $\log p_\theta$ or the score function s_θ via variational inference, for applications one will need to sample from the corresponding distribution $p_\theta(x) dx$. In all methods which facilitate sampling, one defines a stochastic dynamical system such that its states at various time steps are approximate samples from the distribution $p_\theta(x) dx$ of interest. Below we review several prominent samplings methods for sampling from $p_\theta(x) dx$; in our exposition we omit the θ subscript since it will not play a role in the sampling methods.

Markov Chain Monte Carlo (MCMC). A fundamental method, the *Metropolis-Hastings algorithm*, can be used when one has access to the log probability density $\log p(x)$. The algorithm is parameterized by the choice of *proposal distribution* $Q(x|y)$. The method is slightly simplified when $Q(x|y) = Q(y|x)$, namely the *symmetric* case. In this case, one chooses an initial point x_0 from some background distribution, and then runs the following algorithm:

```

initialize  $x_0, i = 0$ 
if  $i < i_{\max}$  do
  sample  $x' \leftarrow Q(x'|x_i)$ 
  set  $x_{i+1} = x'$ 
  with probability  $\alpha = \exp(f(x') - f(x_i)) = p(x')/p(x_i)$ , replace  $i \rightarrow i + 1$ 
end if
return  $x_{i_{\max}}$ 

```

When Q is not symmetric, one simply modifies the formula for the *acceptance ratio* α slightly. For large i_{\max} , under very general conditions [43] the distribution of the variable $x_{i_{\max}}$ will be very

close to $p(x) dx$. For example, choosing $Q(x|y)$ to be a Gaussian centered at y guarantees convergence independent of $p(x)$; here, one can think of the variance of the Gaussian as a “temperature” parameter.

However, while convergence occurs for most choices of a proposal distribution, the *rate of convergence* (how large i_{\max} needs to be for $x_{i_{\max}}$ to be close in distribution to a sample from $p(x)$) and the *decorrelation rate* (a measurement of how large k needs to be so that x_i and x_{i+k} are independently distributed) are both highly sensitive to the choice of proposal distribution Q , and both grow quickly with the dimension and “complexity” of the target distribution p . Both of these quantities affect the accuracy and variance of estimates of expectation values $\mathbb{E}_{x \sim p(x)}[g(x)]$ produced via sampling. Thus, the challenge is to adapt the choice of proposal distribution Q to p in order to make sampling fast enough to be practical.

Score function-based sampling. It is possible to sample from a distribution $p(x) dx$ without having access to $\log p(x)$, but only relying on the score function $s(x)$. The basic method in this class is (*unadjusted*) *Langevin sampling*: one initializes x_0 from an arbitrary initial distribution, and then forms the stochastic process

$$x_{i+1} = x_i + \epsilon \nabla \log p(x) + \sqrt{2\epsilon} \delta w_i \quad (2.11)$$

where $\delta w_i \sim \mathcal{N}(0, \mathbf{1})$ for each i , and uses the fact that under fairly weak regularity conditions $x_i \sim p(x)$ for large i .

If there is access to $\log p(x)$, one can perform *adjusted Langevin sampling* which can be viewed as a hybrid between (unadjusted) Langevin sampling and Metropolis-Hastings. The basic idea is that we can use the stochastic process (2.11) to define a proposal distribution $Q(x|y)$ to be used in the Metropolis-Hastings algorithm. This proposal distribution is useful because it is tailored to the probability distribution $p(x) dx$ of interest. In particular, we define a $Q(x|y)$ via the following sampling algorithm:

```

initialize  $x_0 = y$ 
for  $i = 0, 1, \dots, N$  do
  sample  $\delta w_i \leftarrow \mathcal{N}(0, \mathbf{1})$ 
  set  $x_{i+1} = x_i + \epsilon \nabla \log p(x) + \sqrt{2\epsilon} \delta w_i$ 
end for
set  $x = x_N$ 
return  $x$ 

```

Above, ϵ has to be chosen so that it is sufficiently small, and N has to be chosen so that it is sufficiently large.

Hamiltonian Monte Carlo and related methods. The Langevin sampling algorithm is based on the Euler-Murayama discretization of *overdamped Langevin dynamics*, which is the continuous-time stochastic process defined by the stochastic differential equation (SDE)

$$dx = s(x) dt + \sqrt{2} dW_t. \quad (2.12)$$

Thinking of the log-probability $\log p(x)$ as the negative of a *potential function* $U(x)$, overdamped Langevin dynamics is simply the dynamics of a particle undergoing Brownian diffusion in the potential field U . In Brownian diffusion the particle moves diffusively as opposed to inertially, under the potential U . One can modify the dynamics to incorporate inertial motion under the potential U ; the chaotic mixing properties of such dynamics should heuristically improve the convergence rate of the corresponding stochastic process to the equilibrium distribution $p(x) dx \propto e^{-U(x)} dx$ [45], where in this context the score is $s(x) = -\nabla U(x)$. Such a stochastic process is given by the second-order Langevin dynamics:

$$dx = -v dt, \quad dv = (s(x) - \gamma v) dt + \sqrt{2\gamma} dW_t. \quad (2.13)$$

Here γ is a friction parameter, and one formally recovers overdamped Langevin dynamics by taking $\gamma \rightarrow \infty$, where v changes much faster than x and thus equilibrates to a normal distribution centered at $s(x)$ instantaneously. In the above equation, the $-\gamma v dt$ term is a *friction* term, the $\sqrt{2\gamma} dW_t$ term is a *noise* term, and the $s(x) dt = -\nabla U(x) dt$ term corresponds to noiseless *Hamiltonian dynamics* in the potential U . The ratio between the noise term and the friction term in (2.13) is precisely chosen such that the stable distribution $\tilde{p}(x, v) dx dv$ of the dynamics marginalizes to $p(x) dx$; if the friction term is dropped, this marginalization property no longer holds [46].

Thus, given only access to the score function $s(x)$, one can sample from $p(x) dx$ via *underdamped Langevin sampling* by initializing x_0 , choosing v_0 to be normally distributed around zero, evolving (x_0, v_0) via a discretization of the SDE (2.13) to (x_N, v_N) , and then using x_N as an approximate sample. If one has access the log-probability density, one can, as in adjusted Langevin sampling, instantiate a hybrid method with Metropolis-Hastings; this is called the *adjusted underdamped Langevin sampler*. In this hybrid method, the Metropolis acceptance step renders it unnecessary for the stationary density of the SDE (2.13) to marginalize to $p(x) dx$. Indeed, in this setting one can simply *drop* the friction term (or even both the friction term and the noise term) from the adjusted underdamped Langevin sampler while keeping the Metropolis acceptance step; this gives rise to the *Hamiltonian Monte-Carlo* method which is a mainstay of physical and statistical simulations [12, 45]. One can incorporate Riemannian metrics and variable covariance matrices for the noise in attempts to improve convergence of the resulting samplers [47], and there is an evolving and complex theoretical literature comparing variations of these methods (see e.g. [48–50]).

Bridge sampling. There is a class of methods which, instead of using more complex dynamical systems to define improved proposal distributions $Q(x|y)$, run a series of parallel sampling chains for a collection of different distributions, coupling the chains in some way such that the overall mixing times may be improved. This class of methods includes parallel tempering or replica exchange [10,

51], and is connected to bridge sampling [52]. The philosophy behind these methods is as follows. Consider a sequence or *bridge* of distributions $\mu_0, \mu_1, \dots, \mu_T$ where $\mu = \mu_0 = p(x) dx$ and $\mu_t = p_t(x) dx$. Here μ is the distribution that we want to sample from, and μ_T is a distribution that is easy to sample from (e.g. a unimodal Gaussian), say using an MCMC sampler which converges rapidly and can quickly produce uncorrelated samples. The distributions μ_i should in some sense interpolate between $\mu = \mu_0$ and μ_T . The idea, then, is to use samples from μ_T to seed a proposal distribution to sample from μ_{T-1} , and then to use ensuing samples from μ_{T-1} to seed a proposal distribution to sample from μ_{T-2} , and so on down to $\mu = \mu_0$.

Slightly more explicitly, if $x^{(t)}$ is approximately a sample from μ_t and if $Q_{t-1}(x|y)$ is a proposal distribution for μ_{t-1} , then we can use $Q_{t-1}(x|x^{(t)})$ as a proposal distribution in MCMC to sample from μ_{t-1} . This kind of procedure will tend to produce rapidly mixing proposal distributions, insofar as $\mu_{t-1} \approx \mu_t$ for all t . More generally, in a bridge sampling scheme, one runs samplers for all the μ_t in parallel, transferring samples between the μ_t according to some appropriate rules such that the stationary distributions of the joint chain marginalizes to $\mu = \mu_0$ upon forgetting samples from the auxiliary distributions in the bridge. The method is improved by choosing an appropriate bridge (i.e. judicious choices of the μ_t 's) and appropriate proposal distributions $Q_t(x|y)$ (which in this context are called *exchange* distributions) which are adapted to the features of the μ_t 's. It can be helpful to choose μ_t to have physical meaning: if μ_0 is the distribution of states of a protein at a low temperature, then one can choose the μ_t to be distributions over states of the protein at progressively higher temperatures as t is increased, so that the whole Markov process simulates denaturation and folding of the protein upon temperature cycling [53].

3 Review of mathematical aspects of Diffusion Models

In this section, we review the fundamental ideas underlying a class of score-function-based hierarchical generative model called *diffusion models*, which have recently become the state-of-the-art models for image generation, powering commercially available tools like Midjourney and Stable Diffusion [54]. Some useful reviews include [25, 55, 56].

3.1 Score-based generative modeling

3.1.1 General setup

As mentioned in the previous section, in score-based generative modeling, we attempt to approximate the score $s(x)$ of the true distribution $p(x) dx$ via an estimated score $s_\theta(x)$, which suffices to approximately sample from $p(x)$ via unadjusted Langevin sampling. The natural objective function to use in this case is the ℓ^2 error of the score function under the data, namely

$$\frac{1}{2} \mathbb{E}_{x \sim p(x)} [|s_\theta(x) - \nabla \log p(x)|^2] \tag{3.1}$$

which is just (one half of) the Fisher divergence $D_F(p|p_\theta)$ from (2.5). However, one does not know the true value of $\log p(x)$ in many practical settings. Following Hyvärinen [57], it is prudent to integrate by parts and subtract a θ -independent constant to arrive at the equivalent objective function

$$\mathbb{E}_{x \sim p(x)} \left[\nabla \cdot s_\theta(x) + \frac{1}{2} |s_\theta(x)|^2 \right] \approx \frac{1}{N} \sum_{i=1}^N \left(\nabla \cdot s_\theta(x_i) + \frac{1}{2} |s_\theta(x_i)|^2 \right), \quad (3.2)$$

where x_i , $i = 1, \dots, N$ are samples from $p(x)$ such as e.g. natural images or microscopic configurations of the spins of a magnet. Importantly, (3.2) manifestly does not require explicit knowledge of $\log p(x)$.

3.1.2 Introducing a noise scale

We would like to ameliorate the common issue of $p(x)$ being concentrated along a submanifold of the data space \mathbb{R}^D . A useful approach is to smooth out the probability density $p(x)$ by noising it, forming the new probability density

$$p_\sigma(\tilde{x}) = \int dx p(x) \frac{e^{-\frac{1}{2\sigma^2}|x-\tilde{x}|^2}}{(2\pi\sigma^2)^{D/2}}. \quad (3.3)$$

A convenient way to sample from $p_\sigma(\tilde{x})$ is to take empirical samples $x_i \sim p(x)$ and produce synthetic samples $\tilde{x}_i \sim x_i + \sigma z_i$, where $z_i \sim \mathcal{N}(0, 1)$. One can now try to estimate the score of p_σ for some small value of σ , and use samples from p_σ as an approximation to samples from p . Score matching for p_σ is easier [25, Figure 1], and the corresponding smoothing p_σ^{data} of the empirical data distribution $p^{\text{data}}(x) = \frac{1}{N} \sum_{i=1}^N \delta(x - x_i)$ has a natural gradient estimator. In particular, suppose we parameterize the score function via

$$s_\theta(x) = \frac{1}{\sigma^2} (W^T \cdot \text{sigmoid}(W \cdot x + b) + c - x), \quad \text{where } \theta = (W, b, c) \quad (3.4)$$

and where the dimensions of W are $D' \times D$, the dimensions of b are $D' \times 1$, and the dimensions of c are $D \times 1$. We suppose that $D' < D$. Then minimizing (3.1) using $\tilde{x} \sim p_\sigma^{\text{data}}(\tilde{x})$ in place of $x \sim p(x)$, and also replacing p with p_σ inside the expectation, i.e. minimizing

$$\frac{1}{2} \mathbb{E}_{\tilde{x} \sim p_\sigma^{\text{data}}(\tilde{x})} [|s_\theta(\tilde{x}) - \nabla \log p_\sigma(\tilde{x})|^2], \quad (3.5)$$

turns out to be equivalent (in the sense that the gradients agree up to a global rescaling) to minimizing the quantity [24]

$$\mathbb{E}_{x \sim p(x)} \mathbb{E}_{\substack{\tilde{x} = x + \sigma z \\ z \sim \mathcal{N}(0, \sigma I)}} [|W^T \cdot \text{sigmoid}(W \cdot \tilde{x} + b) + c - x|^2]. \quad (3.6)$$

It is useful to slightly reorganize the above quantity in order to better interpret it. Let us define the nonlinear map

$$\mathcal{E}(\tilde{x}) := \text{sigmoid}(W \cdot \tilde{x} + b) \quad (3.7)$$

which implements an *encoding* of the noised sample into a lower-dimensional (i.e. D' -dimensional) representation, and further define the nonlinear map

$$\mathcal{D}(y) = W^T \cdot y + c \quad (3.8)$$

which will serve as a *decoding* from the lower-dimension representation back to the higher-dimensional (i.e. D -dimensional) one. Then we can rewrite (3.6) as

$$\mathbb{E}_{x \sim p(x)} \mathbb{E}_{\substack{\tilde{x} = x + \sigma z \\ z \sim \mathcal{N}(0, \sigma I)}} \left[|\mathcal{D} \circ \mathcal{E}(\tilde{x}) - x|^2 \right]. \quad (3.9)$$

The above is readily interpreted as the objective function for a *denoising autoencoder*: the composition $\mathcal{D} \circ \mathcal{E}$ is a 1-layer neural network that attempts to *reconstruct* the original sample x from the noised sample \tilde{x} by encoding and subsequently decoding it through a lower-dimensional space.

More broadly, (3.6) for general $s_\theta(x)$ is equivalent to the *denoising score function objective* [24]

$$\mathbb{E}_{x \sim p(x)} \mathbb{E}_{\substack{\tilde{x} = x + \sigma z \\ z \sim \mathcal{N}(0, \sigma I)}} \left[\left| s_\theta(\tilde{x}) - \nabla_{\tilde{x}} \log \left(\frac{e^{-\frac{1}{2\sigma^2}|x-\tilde{x}|^2}}{(2\pi\sigma^2)^{D/2}} \right) \right|^2 \right] = \mathbb{E}_{x \sim p(x)} \mathbb{E}_{\substack{\tilde{x} = x + \sigma z \\ z \sim \mathcal{N}(0, \sigma I)}} \left[\left| s_\theta(\tilde{x}) - \frac{1}{\sigma^2}(x - \tilde{x}) \right|^2 \right]. \quad (3.10)$$

We see that optimizing over θ pressures $s_\theta(\tilde{x})$ to *point in the direction of the denoising of \tilde{x}* , namely the vector $x - \tilde{x}$.

3.1.3 Introducing multiple noise scales

So far we have considered smoothing out $p(x)$ by introducing a noise scale σ . Here we generalize that approach by taking inspiration from bridge sampling. Specifically we introduce a series of noise scales

$$0 < \sigma_1 < \dots < \sigma_T \quad (3.11)$$

and corresponding noised distributions

$$p_{\sigma_i}(\tilde{x}) := \int dx p(x) \frac{e^{-\frac{1}{2\sigma_i^2}|x-\tilde{x}|^2}}{(2\pi\sigma_i^2)^{D/2}}. \quad (3.12)$$

If $s_\theta^{(i)}(x)$ is a score function associated with each p_{σ_i} , we can optimize a weighted sum of score matching losses

$$\sum_{i=1}^L \lambda_i \mathbb{E}_{\tilde{x} \sim p_{\sigma_i}^{\text{data}}(\tilde{x})} \left[\left| \nabla_x \log p_{\sigma_i}(\tilde{x}) - s_\theta^{(i)}(\tilde{x}) \right|^2 \right] \quad (3.13)$$

where the $\lambda_i > 0$ are weighting parameters (usually taken to be $\lambda_i = \sigma_i^2$). The above equation is a multiscale generalization of (3.5). Crucially, taking inspiration from (3.10), one should parameterize $s_\theta^{(i)}(x)$ via a *deep convolutional neural network*; the score functions should be a composition of functions which attempt to find denoisings of noised samples. In practice, one parameterizes the

$s_\theta^{(i)}(x)$ via a U-Net with skip connections [25, 58], a neural network architecture which has proven successful in image modeling. Conveniently, one can sample from $p_{\sigma_1}(x)$ via annealed Langevin sampling: one samples from a uniform distribution, which is a good approximation to p_{σ_T} for large σ_T , then takes several updates as in (2.11) with the score function $s_\theta^{(T)}(x)$, then with $s_\theta^{(T-1)}(x)$, and so on until $s_\theta^{(1)}(x)$. This adapted bridge sampling method then naturally resolves many problems associated with the slow Langevin sampling of the (typically multimodal and highly oscillatory) distribution $p(x)$ with score function $\nabla \log p(x) \approx s_\theta^{(1)}(x)$.

There is a related class of models, referred to as *denoising diffusion models*, which come from a physics-inspired framework [23] based on minimizing a KL divergence (2.4) instead of a score matching objective. In this setting, given an initial sample $x_0 \sim p(x_0)$, one defines noised variables

$$x_{i+1} = \sqrt{1 - \beta_i} x_i + \sqrt{\beta_i} z_i, \quad z_i \sim \mathcal{N}(0, \mathbf{1}). \quad (3.14)$$

The above is a discrete-time Ornstein-Uhlenbeck process such that x_T for large T is approximately distributed according to a mean-zero Gaussian with identity covariance. (This is in contrast with the noising process (3.3), which converges to a uniform distribution.) Slightly overloading notation (albeit in a standard manner), we let $p(x_{i+1}|x_i)$ denote the probability density of x_{i+1} given x_i . We still reserve $p(x_0)$ for denoting the probability density we wish to model. In the above setting of (3.14), $p(x_{i+1}|x_i)$ is a Gaussian in x_{i+1} . It turns out that if β_i is small then $p(x_i|x_{i+1})$ is approximately Gaussian in x_i ; this is the discrete-time analog of the fact that there exists a reverse SDE for every given SDE, which we will recall in more detail shortly. Writing $p(x_0)$ as

$$p(x_0) = \int dx_1 \cdots dx_{T-1} dx_T p(x_0|x_1)p(x_1|x_2) \cdots p(x_{T-1}|x_T) p(x_T) \quad (3.15)$$

where $p(x_T)$ denotes the probability associated with the random variable x_T , we observe that $p(x_T)$ is approximately Gaussian (in x_T) and each of the $p(x_i|x_{i+1})$ is approximately Gaussian (in x_i). This suggests that we can approximate $p(x_0)$ by

$$p_\theta(x_0) = \int dx_1 \cdots dx_{T-1} dx_T p_\theta(x_0|x_1)p_\theta(x_1|x_2) \cdots p_\theta(x_{T-1}|x_T) p_\theta(x_T) \quad (3.16)$$

where we have chosen the parameterizations in the following manner. We let $p_\theta(x_T) = \frac{1}{(2\pi)^{D/2}} e^{-\frac{1}{2}|x_T|^2}$ so that it is θ -independent. We further let $p_\theta(x_i|x_{i+1})$ be a Gaussian in x_i whose mean and variance depend, in a prescribed way, on x_{i+1} and the parameters θ . This set of $p_\theta(x_i|x_{i+1})$'s comprise a useful variational family, as discussed in Section 2.3. If we have judiciously chosen the form of the $p_\theta(x_i|x_{i+1})$'s and optimized for the parameters θ , then an approximate sampling algorithm for $p_\theta(x_0) \approx p(x_0)$ is as follows:

```

sample  $x_T \leftarrow \mathcal{N}(0, \mathbf{1}) = p_\theta(x_T)$ 
for  $i = T, T - 1, \dots, 1$  do
    sample  $x_{i-1} \leftarrow p_\theta(x_{i-1}|x_i)$ 
end for
return  $x_0$ 

```

The above is nice since it only entails sampling from Gaussian distributions at every step.

Following the ideas of variational inference, one recalls (as in Section 2.2) that minimizing $\text{KL}(p(x_0)|p_\theta(x_0))$ agrees with maximizing $\mathbb{E}_{x_0 \sim p(x_0)}[\log p_\theta(x_0)]$, which by (2.10) has the lower bound (see [23])

$$\begin{aligned} \mathbb{E}_{x_0 \sim p(x_0)}[\log p_\theta(x_0)] &\geq \mathbb{E}_{x_0, x_1, \dots, x_T \sim p(x_0, x_1, \dots, x_T)} \left[\log \frac{p_\theta(x_0, x_1, \dots, x_T)}{p(x_1, \dots, x_T|x_0)} \right] \\ &= \mathbb{E}_{x_{T-1}, x_T \sim p(x_{T-1}, x_T)} \left[\log p_\theta(x_T) + \log \frac{p_\theta(x_{T-1}|x_T)}{p(x_T|x_{T-1})} \right] \end{aligned} \quad (3.17)$$

where $p(x_0, x_1, \dots, x_T)$ denotes the joint distribution over all the x_i , and $p(x_{T-1}, x_T)$ denotes the joint distribution over only x_{T-1} and x_T . Via a rearrangement, the quantity in the second line can be rewritten further as

$$\mathbb{E}_{x_0, x_1, \dots, x_T \sim p(x_0, x_1, \dots, x_T)} \left[\text{KL}(p(x_T|x_0)|p_\theta(x_T)) + \sum_{i=2}^T \text{KL}(p(x_{i-1}|x_i, x_0)|p_\theta(x_{i-1}|x_i)) - \log p_\theta(x_0|x_1) \right] \quad (3.18)$$

where here $p(x_{i-1}|x_i, x_0)$ denotes the probability of x_{i-1} conditioned on fixed x_i and x_0 . A useful feature of the above equation is that each of the terms (except the first, which disappears upon taking gradients with respect to θ since $p_\theta(x_T)$ is θ -independent) can be computed analytically, since all the quantities involved are Gaussians which can be explicitly computed. For this purpose is it convenient to use the identity $p(x_{i-1}|x_i, x_0) = \frac{p(x_{i-1}, x_i|x_0)}{\int dx_{i-1} p(x_{i-1}, x_i|x_0)}$, where $p(x_{i-1}, x_i|x_0)$ is an explicitly computable Gaussian using (3.14).

As an explicit example of how to parameterize the $p_\theta(x_{i-1}|x_i)$'s, consider

$$p_\theta(x_{i-1}|x_i) = \frac{1}{(2\pi\beta_i^2)^{D/2}} \exp\left(-\frac{1}{2\beta_i^2} \left|x_{i-1} - \mu_\theta^{(i)}(x_i)\right|^2\right) \quad (3.19)$$

where $\mu_\theta^{(i)}(x_i)$ is given by

$$\mu_\theta^{(i)}(x_i) = \frac{1}{\sqrt{1-\beta_i}} \left(x_i - \epsilon_\theta^{(i)}(x_i)\right) \prod_{j=1}^{i-1} \frac{1}{1-\beta_j}. \quad (3.20)$$

Here $\epsilon_\theta^{(i)}(x_i)$ is parameterized by a suitable neural network. Then optimizing over (3.18), one recovers a particular weighted sum of score matching objectives from a KL-minimization framework. This method led to a practically successful implementation (DDPG) of this paradigm in [16].

3.2 Taking the continuum limit

From the previous discussion we have seen that both score matching and variational inference involve the choice of a process which adds noise to a sample (see (3.3) and (3.14), respectively), as well as the learning of a corresponding denoising process (parameterized by s_θ and ϵ_θ , respectively).

In each case, the noising process occurs in discrete time steps. As such, it is natural to take the continuum limit in time and consider forwards SDEs of the form

$$dx = b(x, t) dt + \sigma(x, t) dW_t \quad (3.21)$$

where we will use the *Itô* formulation.

For example, rewriting the noise process (3.3) as

$$x_i = x_{i-1} + \sqrt{\sigma_i^2 - \sigma_{i-1}^2} z_{i-1}, \quad (3.22)$$

where we have set $\sigma_0 = 0$, we see that the continuum limit of this process (where $\sigma_i = \sigma(i/N)$) as $N \rightarrow \infty$ for some function $\sigma : [0, 1] \rightarrow \mathbb{R}$)

$$dx = \sqrt{\frac{d\sigma(t)^2}{dt}} dW_t. \quad (3.23)$$

On the other hand, the noising process (3.14) has continuum limit²

$$dx = -\frac{1}{2}\beta(t)x dt + \sqrt{\beta(t)} dW_t. \quad (3.24)$$

If we set $x(0)$ to be distributed as $p(x) =: p_{t=0}(x)$, and the process $x(t)$ satisfies (3.21), then $x(T)$ will be distributed according to some distribution $p_{t=T}(x)$, and more generally the marginal distribution of $x(t)$ will be $p_t(x)$ for $t \in [0, 1]$. One can also write down a natural reverse SDE with time parameter $\tilde{t} = T - t$, which flows $p_{t=T}(x)$ to $p_{t=0}(x)$, namely [59]

$$dx = \left(-b(x, \tilde{t}) + \nabla \cdot (\sigma(x, \tilde{t}) \cdot \sigma^T(x, \tilde{t})) + \sigma(x, \tilde{t}) \cdot \sigma^T(x, \tilde{t}) \nabla \log p_{\tilde{t}}(x) \right) d\tilde{t} + \sigma(x, \tilde{t}) dW_{\tilde{t}}. \quad (3.25)$$

What we mean is that if $x(\tilde{t})$ is a process with $x(0)$ distributed as $p_{\tilde{t}=0}(x) = p_{t=T}(x)$, and x satisfies (3.25), then $x(T)$ will be distributed according to $p_{\tilde{t}=T}(x) = p_{t=0}(x)$, and more generally $x(\tilde{t})$ will be distributed according to $p_{\tilde{t}=T-\tilde{t}}(x)$. Notice that the above process involves the score function $\nabla \log p_t(x)$. If $b(x, t)$ and $\sigma(x, t) = g(t)$ are scalars, then (3.25) simplifies to

$$dx = \left(-b(x, \tilde{t}) + g(\tilde{t})^2 \nabla \log p_{\tilde{t}}(x) \right) d\tilde{t} + g(\tilde{t}) dW_{\tilde{t}}. \quad (3.26)$$

With the above notations at hand, we are prepared to take the continuous-time limit of the objective function (3.13), giving us

$$\mathbb{E}_{t \sim [0, T]} \mathbb{E}_{x \sim p_t^{\text{data}}(x)} \left[\lambda(t) |\nabla \log p_t(x) - s_\theta(x, t)|^2 \right] \quad (3.27)$$

where $\lambda(t) \geq 0$ is a weighting function. One can optimize this objective with respect to θ by discretizing time, simulating the forwards process via a discretization, and then estimating (3.27) using the Hyärvinen loss (3.2). Then, having found appropriate parameters θ , one can sample

²Here, this continuum limit is taken in the sense that $\beta_i = (1/N)\beta(i/N)$ for some function $\beta : [0, 1] \rightarrow \mathbb{R}$, as $N \rightarrow \infty$. One shows that the continuum limit is (3.24) using the Taylor expansion of $\sqrt{1-\epsilon}$ in ϵ .

from $p(x) \approx p_{\theta, t=0}(x)$ by simulating the flow of the reverse SDE (3.25). The analog of the identity in (3.10) comes from the fact that we can rewrite (3.27), up to a θ -independent constant, as [55]

$$\mathbb{E}_{t \sim [0, T]} [\lambda(t) \mathbb{E}_{x \sim p(x)} \mathbb{E}_{x' \sim p_{t,0}(x'|x)} [\nabla_{x'} \log p_{t,0}(x'|x) - s_{\theta}(x', t)]] \quad (3.28)$$

where $p_{t,0}(x'|x)$ denotes the probability density that the stochastic process produces x' at time t given that it started at x at time 0. The kernel $p_{t,0}(x'|x)$ is called the time- t *transition kernel* of the noising process (3.21), and in many cases it can be computed analytically. In such cases, using (3.28) as the objective has computational advantages, for instance one can make use of an explicit formula for the kernel $p_{t,0}(x'|x)$ as opposed to solving the SDE by other means.

We can provide a variational bound for $\text{KL}(p_{t=0} | p_{\theta, t=0})$ in terms of (3.27) (or equivalently (3.28)) via an application of the Girsanov theorem [55, 60]. For example, if $b(x, t) = b(t)$ and $\sigma(x, t) = g(t)$ are scalars, then letting $\lambda(t) = \frac{1}{2} g(t)^2$ in (3.27) we have

$$\text{KL}(p_{t=0} | p_{\theta, t=0}) \leq \frac{1}{2} \mathbb{E}_{t \sim [0, T]} [g(t)^2 \mathbb{E}_{x \sim p_t^{\text{data}}(x)} [|\nabla \log p_t(x) - s_{\theta}(x, t)|^2]] + \text{KL}(p_{t=T} | p_{\theta, t=T}). \quad (3.29)$$

where we recall again that $p_{\theta, t=T} = p_{t=\infty}$ is θ -independent.

There is also an equivalent reverse *ODE* which flows p_T to p_0 just as (3.25) does.³ In particular, the equation is

$$\frac{dx}{d\tilde{t}} = -b(x, \tilde{t}) + \frac{1}{2} (\nabla \cdot (\sigma(x, \tilde{t}) \cdot \sigma^T(x, \tilde{t})) + \sigma(x, \tilde{t}) \cdot \sigma^T(x, \tilde{t}) \nabla \log p_{\tilde{t}}(x)) \quad (3.30)$$

where σ^T denotes the transpose of σ .

One now wishes to optimize an approximation to (3.25) over choices of $s_{\theta}(x, t)$ as a proxy for $\nabla \log p_t(x)$. To do this, one discretizes time and then solves (3.21) with an SDE solver, and then backpropagates gradients of the time-discretized loss (3.27) into θ . The backwards SDE/ODE is then implicitly used during training in backpropagation or the adjoint method [61], when computing the gradients of the loss with respect to θ . One can then draw samples from p_0 by drawing samples from p_T , and then solve a numerical discretization of the backwards SDE (3.25) or ODE (3.30) from $\tilde{t} = 0$ to $\tilde{t} = T$. Given that p_T (for large T) usually approximately has an explicit log-likelihood (e.g. p_T is approximately Gaussian), one can get tractable estimates for $\log p_0$ (and its θ -gradients) by using the reverse ODE as well as a formula for the instantaneous change in $\log p_t$ along with the Hutchingson trace estimator [62]. The practicalities of how one designs and trains the neural networks parameterizing the score function, as well as the methods for tuning the SDE, are outside the scope of this paper, but we provide some initial pointers to the literature for interested readers in Appendix A.

³The form of (3.30) follows from the solution to the continuity equation for p_{T-t}

4 Lattice field theory and the renormalization group

4.1 Lattice discretization

Recall that in statistical field theory, we study infinite-dimensional probability distributions over spaces of fields. Formally, these distributions have their log-probability specified by the integral of a local quantity over the domain of the field. In scalar ϕ^4 theory in d dimensions, for which we previously considered the case $d = 2$, the fields are functions $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$, and the probability measure is

$$\frac{1}{Z} e^{-S[\phi]} \mathcal{D}\phi = \frac{1}{Z} \exp \left\{ - \int_{\mathbb{R}^d} dy \left(\frac{1}{2} \nabla\phi \cdot \nabla\phi + \frac{1}{2} m^2 \phi^2 + \frac{\lambda}{4!} \phi^4 \right) \right\} \mathcal{D}\phi, \quad (4.1)$$

where $S[\phi]$ is called the *action* of the field theory, $\mathcal{D}\phi$ is the formal volume measure on the space of fields, and Z is a normalization constant. To make sense of such an infinite-dimensional distribution, which is not obviously well-defined, it is prudent to construct a finite-dimensional analog by fixing a lattice spacing ϵ and a lattice size N , which we take to be even. Then we define lattice points $\epsilon \mathbf{n}$ for $\mathbf{n} \in \{-N/2 + 1, \dots, N/2\}^d$. We will denote $\{-N/2 + 1, \dots, N/2\}$ (modulo N) by \mathbb{Z}_N , and accordingly we write $\mathbf{n} \in \mathbb{Z}_N^d$. One then defines real variables $\phi(\mathbf{n})$ which are meant to be samples of $\phi(y)$ at $y = \epsilon \mathbf{n}$. The derivative in (4.1) is discretized via

$$\partial_i \phi(y) \longrightarrow \frac{1}{\epsilon} (\phi(\mathbf{n} + \mathbf{e}_i) - \phi(\mathbf{n})) \quad (4.2)$$

where we define the vector

$$\mathbf{e}_i := (e_i^1, e_i^2, \dots, e_i^d), \quad e_i^j := \delta_i^j. \quad (4.3)$$

Here δ_i^j is the Kronecker delta. In this paper we will take the lattice to be periodic (i.e. it becomes a latticization of the d -dimensional torus); one could instead impose other boundary conditions which would modify the derivative at the boundary. The integral $\int_{\mathbb{R}^d} dy$ in (4.1) is approximated by the finite volume integral $\int_{(-N/2+1)\epsilon}^{(N/2)\epsilon} \dots \int_{(-N/2+1)\epsilon}^{(N/2)\epsilon} dy$, and then approximately discretized via

$$\int_{(-N/2+1)\epsilon}^{(N/2)\epsilon} \dots \int_{(-N/2+1)\epsilon}^{(N/2)\epsilon} dy f(y) \longrightarrow \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \epsilon^d f(\epsilon \mathbf{n}). \quad (4.4)$$

The finite-dimensional analog of (4.1) is then

$$\mu_{\epsilon, m, \lambda}[\phi_N] = \frac{1}{Z_{N, \epsilon, a, m, \lambda}} \left(\prod_{\mathbf{n} \in \mathbb{Z}_N^d} d\phi(\mathbf{n}) \right) e^{-S_{\epsilon, a, m, \lambda}[\phi_N]} \quad (4.5)$$

where

$$S_{\epsilon, a, m, \lambda}[\phi_N] = \epsilon^d \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \left[\frac{1}{2} \left\{ \sum_{i=1}^d \frac{a^2}{\epsilon^2} (\phi(\mathbf{n} + \mathbf{e}_i) - \phi(\mathbf{n}))^2 \right\} + \frac{1}{2} m^2 \phi(\mathbf{n})^2 + \frac{\lambda}{4!} \phi(\mathbf{n})^4 \right]. \quad (4.6)$$

We have written out the above discretization in full detail in order to be explicit. To obtain the appropriate continuum limit of the action, we would like $a = 1$; however, it will be useful for us to consider the case of general a for reasons we explain later. Notice that when $\lambda > 0$, the $\phi(\mathbf{n})^4$ term dominates the action $S_{\epsilon,a,m,\lambda}[\phi_N]$ for large values of ϕ , and so the probability distribution in (4.5) is normalizable; accordingly, the normalization constant $Z_{N,\epsilon,a,m,\lambda}$ is finite.

Fixing values of N , ϵ , a , m , and λ , the log-density of (4.5) is explicit, so now one can draw samples from the corresponding density using e.g. Hamiltonian Monte Carlo. This is precisely what is done in lattice field theory calculations in order to study physical systems, for instance to estimate the mass of a proton using QCD (although the relevant analog of (4.1) for QCD problems involves many more complications when discretizing the corresponding infinite-dimensional integral, reviewed in [63]). Much work in lattice field theory involves using the structure of probability distributions like (4.5) in order to build more efficient samplers.

4.2 Renormalization group

4.2.1 The lattice setting

In order to define (4.5), we introduced discretization parameters N and ϵ such that we are formally sampling $\phi(y)$ over a cube of side length $N\epsilon$ on densely-packed lattice points with lattice spacing ϵ . However, these discretization parameters do not appear in the formal expression (4.1). Since the sampling scale ϵ that we introduced was *arbitrary*, we would like for the physical predictions of the model, at least for distances *much larger* than the sampling scale ϵ , to be independent of ϵ . Let us fix $N\epsilon = L$, and write $\mu_{N,m,\lambda}$ for the distribution (4.5); here we have dropped the dependence on ϵ since it is fixed in terms of N by $\epsilon = L/N$. Since the dimensionality of the distribution $\mu_{N,a,m,\lambda}$ varies with N , in order to *make sense* of the independence of the long-range physics on N (and ϵ) we need to be able to *compare* between between the distributions $\mu_{N,a,m,\lambda}$ for different values of N . Moreover, in order for the long-range physics to be independent of N (and ϵ), it may be necessary to vary the parameters a , m and λ with N ; as such, as write $\mu_{N,a_N,m_N,\lambda_N}$ to allow for this possibility.

Since samples ϕ_N from $\mu_{N,a_N,m_N,\lambda_N}[\phi_N]$ are supposed to be lattice approximations to the continuous field $\phi(y)$, it is natural to implement *comparison maps* by defining some “resampling” or “field interpolation” map $f_{NN'}$ for $N > N'$, which (deterministically or stochastically) takes samples ϕ_N from $\mu_{N,a_N,m_N,\lambda_N}[\phi_N]$ to samples $\phi_{N'}$ from $\mu_{N',a_{N'},m_{N'},\lambda_{N'}}[\phi_{N'}]$. We would like for such maps $f_{NN'}$ to be *local*. This means that $[f_{NN'}(\phi_N)](\mathbf{n}')$, which is a field configuration $\phi_{N'}(\mathbf{n}')$ evaluated at the location $\frac{L}{N'} \mathbf{n}' \in \mathbb{R}^d$, should be computed from *nearby lattice sites* at lattice spacing L/N ; in other words, $[f_{NN'}(\phi_N)](\mathbf{n}')$ is primarily a function of $\phi(\mathbf{n})$ for $|\frac{L}{N} \mathbf{n} - \frac{L}{N'} \mathbf{n}'| \leq cL/N$ for some $c \sim O(1)$. For example, on a 2D lattice if we fix $N_t = 2^t$ for $t = 0, 1, 2, \dots$, such that as t increases the lattice sites get dyadically subdivided, then it is natural to choose $[f_{N_{t+1},N_t}(\phi(\mathbf{n}))](\mathbf{n}')$ to be the the average of the values of $\phi(\mathbf{n})$ on the 2×2 square subdivided sites.

The locality of the comparison map enforces that certain long-range properties of correlations

are preserved upon resampling. To be precise, it is useful to define *correlation functions* which are expectation values of a product of fields at different lattice sites as a function of the position:

$$c_{\mu_N, a_N, m_N, \lambda_N}[\phi_N](\mathbf{n}_1, \dots, \mathbf{n}_r) := \mathbb{E}_{\phi_N \sim \mu_N, a_N, m_N, \lambda_N}[\phi_N][\phi(\mathbf{n}_1)\phi(\mathbf{n}_2) \cdots \phi(\mathbf{n}_r)]. \quad (4.7)$$

It is also convenient to define the resampled correlation functions

$$c_{\mu_N, a_N, m_N, \lambda_N}^{f_{NN'}}[\phi_N](\mathbf{n}'_1, \dots, \mathbf{n}'_r) := \mathbb{E}_{\phi_N \sim \mu_N, a_N, m_N, \lambda_N}[\phi_N][f_{NN'}(\phi_N)(\mathbf{n}'_1)[f_{NN'}(\phi_N)(\mathbf{n}'_2) \cdots [f_{NN'}(\phi_N)(\mathbf{n}'_r)]], \quad (4.8)$$

which we can think of as capturing correlations of $\mu_N, a_N, m_N, \lambda_N$ at intermediate to large distance scales; the resampled correlation functions do not capture correlations at short distance scales since the resampled fields $[f_{NN'}(\phi_N)](\mathbf{n}')$ are not strictly local on the lattice. It follows from the usual properties of pushforwards that

$$c_{\mu_N, a_N, m_N, \lambda_N}^{f_{NN'}}(\mathbf{n}'_1, \dots, \mathbf{n}'_r) = c_{(f_{NN'})_* \mu_N, a_N, m_N, \lambda_N}(\mathbf{n}'_1, \dots, \mathbf{n}'_r), \quad (4.9)$$

which means that the intermediate to long-distance properties of correlation functions of $\mu_{a_N, m_N, \lambda_N}[\phi_N]$ are preserved in $(f_{NN'})_* \mu_{a_N, m_N, \lambda_N}[\phi_N]$.

In many circumstances, it is useful to approximate $(f_{NN'})_* \mu_{a_N, m_N, \lambda_N}[\phi_N]$ by $\mu_{N', a_{NN'}, m_{NN'}, \lambda_{NN'}}[\phi_{N'}]$ with judiciously chosen couplings $a_{NN'}$, $m_{NN'}$, and $\lambda_{NN'}$. In particular, we seek parameters $a_{NN'}$, $m_{NN'}$, and $\lambda_{NN'}$ such that

$$c_{(f_{NN'})_* \mu_{a_N, m_N, \lambda_N}[\phi_N]}(\mathbf{n}_1, \dots, \mathbf{n}_r) \approx c_{\mu_{N', a_{NN'}, m_{NN'}, \lambda_{NN'}}[\phi_{N'}]}(\mathbf{n}_1, \dots, \mathbf{n}_r) \quad (4.10)$$

for $|\frac{L}{N'} \mathbf{n}_i - \frac{L}{N'} \mathbf{n}_j|_{\mathbb{Z}_N^d} \gg O(1)$ with $i \neq j$. The notation $|\mathbf{n} - \mathbf{m}|_{\mathbb{Z}_N^d}$ means the (shortest) distance between \mathbf{n} and \mathbf{m} on \mathbb{Z}_N^d viewed as a periodic lattice. This is all to say that we want the intermediate to long-distance correlation functions of $(f_{NN'})_* \mu_{a_N, m_N, \lambda_N}[\phi_N]$ and $\mu_{N', a_{NN'}, m_{NN'}, \lambda_{NN'}}[\phi_{N'}]$ to approximately match, if it is possible to do so.

The above discussion of how the coupling effectively change on account of the flow is so important that it is worth drawing it out more explicitly. Let us write $\mathcal{D}\phi_N := \prod_{\mathbf{n} \in \mathbb{Z}_N^d} d\phi(\mathbf{n})$ for the Euclidean measure over fields sampled at lattice spacing L/N . Then we can write $(f_{NN'})_* \mu_{a_N, m_N, \lambda_N}$ in a similar form as $\mu_{N', a_{NN'}, m_{NN'}, \lambda_{NN'}}$:

$$(f_{NN'})_* \mu_{a_N, m_N, \lambda_N}[\phi_N] = \frac{1}{Z_{N'}} \exp \left(- \left(\frac{L}{N'} \right)^d \sum_{\mathbf{n} \in \mathbb{Z}_{N'}^d} \left[\frac{1}{2} \left\{ \sum_{i=1}^d \frac{a_{NN'}^2}{\left(\frac{L}{N'} \right)^2} (\phi(\mathbf{n} + \mathbf{e}_i) - \phi(\mathbf{n}))^2 \right\} \right. \right. \quad (4.11)$$

$$\left. \left. + \frac{1}{2} m_{NN'}^2 \phi(\mathbf{n})^2 + \lambda_{NN'} \phi(\mathbf{n})^4 \right] + \cdots \right) \mathcal{D}\phi_{N'}$$

where $1/Z_{N'}$ is an appropriate normalization. Moreover the \cdots in the exponential denote terms in the log probability density contain other polynomials in the $\phi(\mathbf{n})$'s (and possibly non-polynomials) which correspond to non-local and higher-order interactions. Insofar as the \cdots terms are small, we

can choose not to include them; then we can more cleanly think of a_N , m_N , and λ_N as flowing to $a_{NN'}$, $m_{NN'}$, and $\lambda_{NN'}$, respectively. In particular, what it *means* for the \dots terms to be small is that (4.10) holds for intermediate to long-distance correlators.

Thus, we can think of the *averaging* process $f_{NN'}$ as generating a *dependence* of the coefficients a_N, m_N, λ_N on N . Alternatively, the requirement (4.10), which expresses that the long-range correlation functions are independent of N , forces us to *tune* the parameters a_N, m_N, λ_N as a function of N such that this physical requirement holds. We conclude that *the coefficients of the model are dependent on the discretization scale of the field*. This dependence of the coefficients of the model on the discretization scale belongs to the framework of the *renormalization group* (or *RG*). A choice of *renormalization group scheme* is essentially a choice of averaging process $f_{NN'}$ for all (N, N') such that $f_{N'N''} \circ f_{NN'} = f_{NN''}$ (the *semigroup property*).

In some circumstances, it is convenient to rescale the ϕ field so that the *kinetic* term in the action, namely $\frac{1}{2} \sum_{i=1}^d \frac{1}{\left(\frac{L}{N}\right)^2} (\phi(\mathbf{n} + \mathbf{e}_i) - \phi(\mathbf{n}))^2$, has a unit coefficient. In a sense, this term sets the size of fluctuations of the field ϕ , and so it may be useful to work in units where such fluctuations are canonically normalized. Concretely, suppose that at scale $\epsilon = L/N$ we have $a_N = 1$. Then performing one RG step to scale $\epsilon' = L/N'$, we may obtain an $a_{NN'}$ which does not equal one. However, implementing the field redefinition $\phi(\mathbf{n}) \rightarrow \phi(\mathbf{n})/a_{NN'}$ and defining $\tilde{m}_{NN'} := m_{NN'}/a_{NN'}$ as well as $\tilde{\lambda}_{NN'} := \lambda_{NN'}/a_{NN'}^4$, (4.11) becomes

$$(f_{NN'})_* \mu_{m_N, \lambda_N}[\phi_N] = \frac{1}{\tilde{Z}_{N'}} \exp \left(- \left(\frac{L}{N'} \right)^d \sum_{\mathbf{n} \in \mathbb{Z}_{N'}^d} \left[\frac{1}{2} \left\{ \sum_{i=1}^d \frac{1}{\left(\frac{L}{N'}\right)^2} (\phi(\mathbf{n} + \mathbf{e}_i) - \phi(\mathbf{n}))^2 \right\} \right. \right. \\ \left. \left. + \frac{1}{2} \tilde{m}_{NN'}^2 \phi(\mathbf{n})^2 + \tilde{\lambda}_{NN'} \phi(\mathbf{n})^4 \right] + \dots \right) \mathcal{D}\phi_{N'} \quad (4.12)$$

where $\tilde{Z}_{N'}$ is the new normalization appropriate for the distribution. The equation above now does have the desired kinetic term, at the cost of modifying the definition of ϕ . Then the analog of (4.10) is

$$c_{(f_{NN'})_* \mu_{N,1,m_N,\lambda_N}[\phi_N]}(\mathbf{n}_1, \dots, \mathbf{n}_r) \approx \frac{1}{a_{NN'}^r} c_{\mu_{N',1,\tilde{m}_{NN'},\tilde{\lambda}_{NN'}}[\phi_{N'}]}(\mathbf{n}_1, \dots, \mathbf{n}_r) \quad (4.13)$$

for $|\frac{L}{N'} \mathbf{n}_i - \frac{L}{N'} \mathbf{n}_j|_{\mathbb{Z}_N^d} \gg O(1)$ with $i \neq j$, where we have accounted for the rescaling of the field ϕ with the factors of $a_{NN'}$. A way to think about (4.13) is that given a $\mu_{N,1,m_N,\lambda_N}[\phi_N]$, we can ask for parameters $a_{NN'}$, $\tilde{m}_{NN'}$, and $\tilde{\lambda}_{NN'}$ such that (4.13) holds for intermediate to long-distance correlators. The formulation preserves the canonical normalization of the kinetic term in the action after each RG step.

To conclude this subsection, we emphasize that the infinite-dimensional distributions (4.1) are ill-defined, and (like all integrals) need to be defined via a limiting procedure. The dependence of the coefficients a_N, m_N, λ_N on the discretization scale N , is not illusory; indeed, it is required for the *definition* of (4.1) as a limit of measures such as (4.5) or (4.11). This phenomenon repeatedly arises in works in mathematical physics on constructive statistical field theory, which aim to make

rigorous mathematical sense of measures like (4.1) [64–66]. In many cases, some of the coefficients a_N, m_N, λ_N are forced to go to zero or to infinity as $N \rightarrow \infty$ in order to get a well-defined limiting distribution.

4.3 A toy analog of the renormalization group

While the procedure described above may seem complicated, a simple of the same phenomenon already arises in the setting of the simplest continuous-time stochastic process, namely Brownian motion. Letting $X(t)$ be a Brownian motion such that $dX(t) = dB_t$, we discretize time into steps of size $\delta t = \epsilon$. This amounts to replacing $X(t)$ with $X^\epsilon(i)$, and setting $X^\epsilon(i+1) = X^\epsilon(i) + \sqrt{\epsilon} Z_i$ where $Z_i \sim \mathcal{N}(0, 1)$. Then the joint distribution over $X^\epsilon(i)$ and $X^\epsilon(i+1)$ becomes

$$\frac{1}{Z} \exp\left(-A_\epsilon X^\epsilon(i)^2 - B_\epsilon X^\epsilon(i+1)^2 - C_\epsilon (X^\epsilon(i) - X^\epsilon(i+1))^2\right), \quad (4.14)$$

which is analogous to the quadratic part of (4.11), where $A_\epsilon, B_\epsilon, C_\epsilon$ are chosen such that

$$\mathbb{E}[X^\epsilon(i)^2] = i \epsilon \quad (4.15)$$

$$\mathbb{E}[X^\epsilon(i+1)^2] = (i+1)\epsilon \quad (4.16)$$

$$\mathbb{E}[(X^\epsilon(i+1) - X^\epsilon(i))^2] = \epsilon. \quad (4.17)$$

We are forced to rescale $A_\epsilon, B_\epsilon, C_\epsilon$ this way with ϵ such that the $\epsilon \rightarrow 0$ limit gives a well-defined continuous stochastic process. Donsker’s theorem [67] shows that we can define the same continuous-time stochastic process $X(t)$ from the $\epsilon \rightarrow 0$ limit of a large family of discrete-time stochastic processes $X^\epsilon(i)$: we only need to require that $X^\epsilon(i+1) - X^\epsilon(i)$ is mean-zero and of variance ϵ , while its Gaussianity is unnecessary. The independence of the *macroscopic* properties of the model from its *microscopic* specification is a simple analog of the *universality* phenomenon connected to RG which is used in statistical field theory to make sense of phase transitions of physical systems, which we briefly overview in Section 4.5.

4.4 Interfacing the renormalization group with modeling and simulation

To make sense of drawing samples from the density (4.1), we have introduced a discretization parameter ϵ so that we can instead sample from discrete analogs like (4.12). In the context of e.g. (4.12), we have argued that given an initial m_N and λ_N (and $a_N = 1$) at scale $\epsilon = L/N$, we can determine the flowed parameters $m_{N'} := \tilde{m}_{NN'}$ and $\lambda_{N'} := \tilde{\lambda}_{NN'}$ (as well as $a_{N'} := a_{NN'}$) at some other fixed value of $\epsilon' = L/N'$ for $N' < N$. This is to say that if we know m_N and λ_N at some initial scale $\epsilon = L/N$, then we can determine a (approximate) flow of the theory to all larger scales; this means that m_N and λ_N at that initial scale serve as initial conditions.

When comparing with a physical system, we would like to use some experimental measurements to set the initial conditions m_N and λ_N of our model. Then the model can be leveraged via simulations (e.g. RG flow and sampling) to make predictions about the experimental system which

can be compared with data. The fixing of the initial conditions is achieved by *scale-setting*: one computes an expectation value of some quantities $\mathcal{F}_i(\phi_N)$ over $\phi_N \sim \mu_{1,m_N,\lambda_N}[\phi_N]$, where the \mathcal{F}_i are experimentally accessible functions of the physical field $\phi(x)$ that can be written in terms of some correlation functions at scale $\epsilon = L/N$. One then sets the values of m_N and λ_N such that the experimentally measured values $\mathcal{F}_{\text{expt},i}[\phi]$ are approximately the predicted values of $\mathcal{F}_i(\phi_N)$ using samples from $\mu_{N,1,m_N,\lambda_N}[\phi_N]$. Only a finite number of experimental measurements are thus needed to set the parameters, and once the parameters are set, other experimental predictions can be made from the simulated model. In practice, because much is known about the structure of the renormalization group for physically relevant models, other RG-based methods can be used to set the scale, such as the method of *gradient flow* [27], which is closely connected to the ideas of this paper.

The quantities m, λ in (4.1) are, in the context of other models, interpreted as a “mass” and “interaction strength”, respectively. Thus, the renormalization group suggests that the experimentally-measured mass of a particle should be dependent on the precision of a measurement used to measure its mass, which in fact occurs experimentally; the same is true of the interaction strength (see e.g. standard texts like [68, 69]).

4.5 Comments on RG fixed points and phase transitions

Here we make comments about a more general class of models. Suppose that our measure has the form

$$\mu_{N,\{\lambda_{i,N}\}}[\phi_N] = \frac{1}{Z_N} \exp\left(-\left(\frac{L}{N}\right)^d \sum_i \lambda_{i,N} M_{i,N}[\phi_N]\right) \mathcal{D}\phi_N \quad (4.18)$$

where the $M_{i,N}[\phi_N]$'s given by

$$M_{i,N}[\phi_N] = \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \prod_{k=0}^{k_{\max,i}} (\Delta^k \phi(\mathbf{n}))^{q_{i,k}} \quad (4.19)$$

where each $q_{i,k} \in \mathbb{Z}_{\geq 0}$ and Δ is the discrete Laplacian operator which acts on fields as

$$\Delta \phi(\mathbf{n}) = \sum_{i=1}^d \frac{1}{\left(\frac{L}{N}\right)^2} (\phi(\mathbf{n} + \mathbf{e}_i) - 2\phi(\mathbf{n}) + \phi(\mathbf{n} - \mathbf{e}_i)) . \quad (4.20)$$

By Δ^k we simply mean applying Δ to a function k times. It will be convenient to define

$$D_i := \sum_{k=0}^{k_{\max,i}} q_{i,k} \quad (4.21)$$

which counts the total number of multiplicative ϕ 's appearing in an $M_{i,N}$. Note that $M_{i,N}[\phi_N]$ depends on N in two ways: through ϕ_N and the sum over $\mathbf{n} \in \mathbb{Z}_N^d$, and through the powers of $\frac{L}{N}$

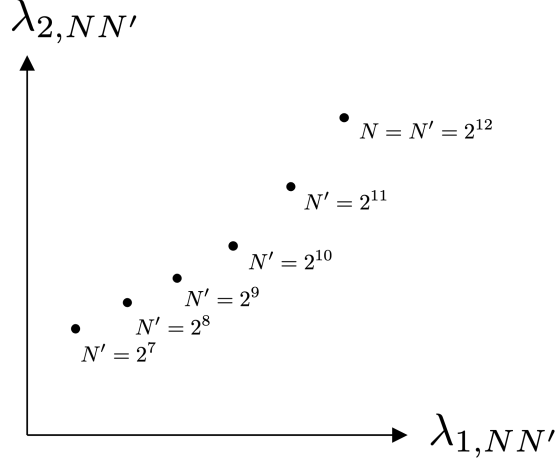


Figure 4: Schematic of the flow of $\lambda_{1,NN'}$ and $\lambda_{2,NN'}$ with varying N' , plotted as $\lambda_{2,NN'}$ versus $\lambda_{1,NN'}$.

coming from the Laplacians. When we write $M_{i,N'}[\phi_{N'}]$, we mean that we are changing all of those dependencies on N to N' .

We can (approximately) flow the couplings to smaller values of N , namely $N' < N$, using an RG flow map. In particular, we define the couplings $\lambda_{i,NN'}$ by

$$(f_{NN'})_* \mu_{N,\{\lambda_{i,N}\}}[\phi_N] \approx \mu_{N',\{\lambda_{i,NN'}\}}[\phi_{N'}] = \frac{1}{Z_{N'}} \exp\left(-\left(\frac{L}{N'}\right)^d \sum_i \lambda_{i,NN'} M_{i,N'}[\phi_{N'}]\right) \mathcal{D}\phi_{N'}, \quad (4.22)$$

where the right-hand side of the \approx should be regarded as the best approximation (in some quantitative way that can be specified) to the left-hand side. Suppose that for L sufficiently large, we plot the trajectories of the couplings $\{\lambda_{i,NN'}\}_i$ as a function of decreasing N' . As we decrease N' (for $1 \ll N' \ll L$), we obtain a diagram as in Figure 4, which should be thought of as a plot of the RG flow of the parameters of the model. There can arise *special points* in this plot which appear to be fixed points of the flow, in the following sense: there exists a constant y such that for $1 \ll N'' < N' \leq N$,

$$\lambda_{i,NN'} \approx \left(\frac{N'}{N''}\right)^{(d-y) \cdot D_i} \lambda_{i,NN''} \quad \text{for all } i. \quad (4.23)$$

Letting $\hat{\lambda}_i := \lambda_{i,NN'}$, the above equations mean that

$$(f_{N'N''})_* \mu_{N',\{\hat{\lambda}_i\}}[\phi_{N'}] \approx \mu_{N'',\{\hat{\lambda}_i\}}\left[\left(\frac{N'}{N''}\right)^{d-y} \phi_{N''}\right]. \quad (4.24)$$

In words, if we perform RG on $\mu_{N',\{\hat{\lambda}_i\}}[\phi_{N'}]$, the distribution has approximately the same couplings, and only the $\phi(\mathbf{n})$ fields are rescaled. Such distributions thus possess a type of scale-invariance, and are called RG fixed points. In nature, *second-order phase transitions* correspond to RG fixed

points. This includes critical points of water, at which the distinction between liquid and gaseous water ceases to exist [70].

The scale-invariance of a critical point, as expressed in (4.24), can manifest in correlation functions. For instance, if the interactions in (4.18) are approximately isotropic with respect to the lattice, then we expect the two point function of the scalar to approximately satisfy

$$c_{\mu_{N'}, \{\hat{\lambda}_i\}, [\phi_{N'}]}(\mathbf{n}_1, \mathbf{n}_2) \approx \frac{C}{|\mathbf{n}_1 - \mathbf{n}_2|_{\mathbb{Z}_{N'}^d}^{2(d-y)}} \quad (4.25)$$

for $1 \ll |\mathbf{n}_1 - \mathbf{n}_2|_{\mathbb{Z}_{N'}^d} \ll N'$, where C is a constant. See e.g. [71] for an explanation of how this kind of power law scaling arises. In any case, one often finds that correlation functions at critical points have a power law behavior.

More broadly, in the limit of large N we can think of the class of distributions (4.18) as comprising the space of *all possible field theories* of a translation-invariant scalar in a fixed number of spatial dimensions. From this point of view, RG flows can be thought of as implementing dynamics on the space of theories: we start with an initial theory, and move through the space of theories by using the RG flow map. Although the dynamics through theory space is contingent on the particular choice of RG transformation f , the *qualitative* features of the RG flow diagram can be independent of the choice of RG scheme. It is a remarkable empirical fact, with a theoretical justification that is outside the scope of this paper (see e.g. [13, 14, 71, 72]), that for many different models like (4.5) RG flows possess *quantitative* features that are largely independent of the choice of RG transformation. This is referred to in the literature as *scheme independence*. For instance, quantities like y in (4.25) associated to an RG fixed point are scheme-independent. Moreover, the general theory of RG flows suggests that we can truncate the sum in the exponential of (4.18) to a small number of terms, where the choice of those terms is contingent on the initial distribution from which we start the RG flow. Such methodologies are central to the study of statistical models of extended physical systems.

4.6 Continuous-time formalism

So far, our discussion of RG on the lattice has considered discrete RG steps, i.e. in which we coarse-grain the lattice from N^d sites to $(N')^d < N^d$ sites. We likened this to a discrete-time dynamical system on the space of lattice theories with a fixed number of sites N^d , where time is a proxy for the number of iterations of the RG flow. However, there is a similar, alternative framework that enables *continuous* RG steps; the induced RG flows can then be viewed as continuous-time dynamical systems on the space of lattice theories. In this context, time is likewise a proxy for the amount that we have flowed or coarse-grained.

To understand how continuous RG is possible on the lattice, we begin with a high-level discussion before delving into particular continuous RG schemes. For simplicity, suppose we again have a lattice scalar field on a d -dimensional lattice of size N^d , which defines lattice points $\epsilon \mathbf{n}$ for $\mathbf{n} \in \mathbb{Z}_N^d$. We further equip the lattice with periodic boundary conditions. We can thus think of the lattice

as living on a d -dimensional torus. In this setting, we can readily take the Fourier transform of the field $\phi(\mathbf{n})$, namely

$$\tilde{\phi}(\mathbf{p}) := \frac{1}{N^d} \sum_{\mathbf{n} \in \mathbb{Z}_N^d} e^{-i \frac{2\pi}{N} \mathbf{p} \cdot \mathbf{n}} \phi(\mathbf{n}) \quad (4.26)$$

where $\mathbf{p} \in \mathbb{Z}_N^d$ (i.e. $\mathbf{p} \in \{-N/2+1, \dots, 0, 1, \dots, N/2\}^d$ (modulo N)). The high-frequency components of $\tilde{\phi}(\mathbf{p})$ correspond to large $|\mathbf{p}|$. For reasons we will discuss later, it will be prudent to consider $\widehat{\mathbf{p}}$ which has components

$$\widehat{p}_i := \frac{2N}{L} \sin\left(\frac{2\pi}{N} \frac{p_i}{2}\right), \quad i = 1, \dots, d, \quad (4.27)$$

so that $|\widehat{\mathbf{p}}|^2 = \sum_{i=1}^d \widehat{p}_i^2$. Using this notation, we note that high-frequency components of $\tilde{\phi}(\mathbf{p})$ also correspond to large $|\widehat{\mathbf{p}}|$.

Thus if we want to smooth out $\phi(\mathbf{n})$ in real-space, then we can equivalently progressively remove the high-frequency components of $\tilde{\phi}(\mathbf{p})$.

4.6.1 A first guess

From the above point of view, an initial guess for a continuous RG map is

$$f_t^{\text{guess}}[\tilde{\phi}](\mathbf{p}) = e^{-|\widehat{\mathbf{p}}|^2 t} \tilde{\phi}(\mathbf{p}), \quad (4.28)$$

for $t \geq 0$. The map $f_t^{\text{guess}}[\tilde{\phi}](\mathbf{p})$ progressively dampens the higher-frequency modes as we increase t . In particular, modes with $|\widehat{\mathbf{p}}| \gg \frac{1}{\sqrt{t}}$ are significantly dampened, and so here $\Lambda_t \approx \frac{1}{\sqrt{t}}$ is our effective *UV cutoff* scale. Note that the lattice scale itself also provides a maximum cutoff, since $|\widehat{\mathbf{p}}|$ can be at most $2\sqrt{d}N/L$.

However, $f_t^{\text{guess}}[\tilde{\phi}](\mathbf{p})$ is in fact *not* a valid RG scheme. This fact was emphasized in [8, 73]; let us provide some of the essential intuition here. Suppose we measure a system such that we can only access momentum modes less than a scale Λ ; we call these modes *IR modes*. Accordingly we have the inability to access momentum modes greater than the scale Λ ; we call these modes *UV modes*. In the microscopic distribution, there are IR modes which are coupled to the UV modes; as such, any reasonable procedure for marginalizing over the UV modes will affect the residual distribution over IR modes which can measure. This can be articulated much more concretely: any physical RG map $f_t[\tilde{\phi}](\mathbf{p})$ should be a mixture of momentum modes when $|\widehat{\mathbf{p}}|$ is near the cutoff scale. For instance, maps of the form⁴

$$f_t[\tilde{\phi}](\mathbf{p}) = \sum_{\mathbf{q}} Q_t(|\widehat{\mathbf{p}} - \widehat{\mathbf{q}}|) e^{-|\widehat{\mathbf{q}}|^2 t} \tilde{\phi}(\mathbf{q}) \quad (4.29)$$

for non-trivial Q_t do mix momentum modes. This requirement of mixing momentum modes is achieved by block-spin RG, but is *not* achieved by the map $f_t^{\text{guess}}[\tilde{\phi}](\mathbf{p})$ in (4.28) which merely rescales each $\tilde{\phi}(\mathbf{p})$ and hence does not mix momentum modes.

⁴More generally, the mixture of momentum modes could be nonlinear.

There is a way to augment the map $f_t^{\text{guess}}[\tilde{\phi}](\mathbf{p})$ so that it can become a valid RG scheme. Suppose that we construct a new map $f_t[\tilde{\phi}](\mathbf{p})$ which is *stochastic*, satisfying

$$\mathbb{E}[f_t[\tilde{\phi}](\mathbf{p})] = e^{-|\hat{\mathbf{p}}|^2 t} \tilde{\phi}(\mathbf{p}), \quad (4.30)$$

where we are averaging over some stochastic noise. If the stochastic map has non-trivial *variance* in the sense that $\tilde{\phi}(\mathbf{p})$ can be stochastically mapped into other nearby momentum modes, then our stochastic map can comprise a valid RG scheme. We pursue this approach below.

4.6.2 A first look at stochastic RG and the Carosso scheme

Let us formulate a stochastic analog of (4.28) along the lines of (4.30), defined implicitly through a systems of stochastic ODEs. Since the expressions are straightforward in position space (i.e. using \mathbf{n} instead of \mathbf{p}), we will work in position space here. In particular, the Carosso scheme [8] considers

$$f_t[\phi](\mathbf{n}) =: \phi_t(\mathbf{n}) \quad (4.31)$$

as the solution to the stochastic differential equation

$$\partial_t \phi_t(\mathbf{n}) = \Delta \phi(\mathbf{n}) + \eta_t(\mathbf{n}), \quad \phi_0(\mathbf{n}) = \phi(\mathbf{n}), \quad (4.32)$$

for all $\mathbf{n} \in \mathbb{Z}_N^d$, where here Δ denotes the discrete Laplacian defined in (4.20) and $\eta_t(\mathbf{n})$ is a Gaussian random field satisfying

$$\mathbb{E}[\eta_t(\mathbf{n})] = 0, \quad \mathbb{E}[\eta_t(\mathbf{n})\eta_s(\mathbf{m})] := \left(\frac{N}{L}\right)^d \delta(t-s) \exp\left(-\Lambda_0^2 \epsilon^2 |\mathbf{n} - \mathbf{m}|_{\mathbb{Z}_N^d}^2\right), \quad (4.33)$$

again for all $\mathbf{n}, \mathbf{m} \in \mathbb{Z}_N^d$. Since $\exp\left(-\Lambda_0^2 \epsilon^2 |\mathbf{n} - \mathbf{m}|_{\mathbb{Z}_N^d}^2\right)$ equals 1 for $\mathbf{n} = \mathbf{m}$ and is approximately 0 for $\mathbf{n} \neq \mathbf{m}$, we can replace the Gaussian by a Kronecker delta $\delta_{\mathbf{n}, \mathbf{m}}$ so that (4.33) becomes

$$\mathbb{E}[\eta_t(\mathbf{n})] = 0, \quad \mathbb{E}[\eta_t(\mathbf{n})\eta_s(\mathbf{m})] := \left(\frac{N}{L}\right)^d \delta(t-s) \delta_{\mathbf{n}, \mathbf{m}}. \quad (4.34)$$

We can think of the Carosso RG flow as smoothing out the scalar field ϕ at progressively larger distance scales as t becomes larger. In momentum space, the Carosso scheme satisfies

$$\mathbb{E}[\tilde{\phi}_t(\mathbf{p})] = e^{-|\hat{\mathbf{p}}|^2 t} \tilde{\phi}(\mathbf{p}), \quad (4.35)$$

but with non-trivial variance in the sense that $\tilde{\phi}(\mathbf{p})$ can be stochastically mapped to other nearby momentum modes. In this sense, the Carosso scheme provides a realization of (4.30).

We emphasize that we can view (4.31) as being a continuous RG map, where the continuum parameter is t . Increasing t means we are further along the RG flow. A key feature of (4.31) is that the lattice remains the same size for any t . That is, unlike our previous RG schemes which change the lattice size, here we are instead affecting the profiles of the fields that can live on a lattice of

fixed size. In fact, our lattice analyses in the previous subsections can be carried over to continuous setting with appropriate modifications.

There are many different kinds of stochastic RG flows, for instance the Polchinski RG flow [34], which we will discuss later. We note that the Carosso RG flow and its cousins are all lattice discretizations of RG flows for continuous fields on \mathbb{R}^d . These RG flows for continuous fields are part of the formalism of the Exact Renormalization Group (ERG), which we review in Appendix C. In the ERG context, the stochastic ODE defined in (4.32) becomes a stochastic PDE. The reason is that in the ERG context the fields are defined for all spatial points $\mathbf{y} \in \mathbb{R}^d$ as opposed to on a finite lattice.

4.7 The concept of Effective Field Theory

4.7.1 General comments

So far, we have explored various manifestations of RG flow in the lattice setting. In particular, we have emphasized the role of *RG fixed points*, which can be viewed as endpoints of RG flows that can possess a type of scale-invariance. But there is another key concept which RG flows enable us to formulate, namely Effective Field Theory.

The basic idea is as follows. Suppose we have some *microscopic theory*, meaning a probability distribution describing degrees of freedom which are not arbitrarily close together in space. Our lattice theories for finite ϵ are an example. In the parlance of physics, such theories are also referred to as *UV theories*, stemming from the terminology of UV light being short-wavelength (at least relative to visible light). We often suppose that the UV theories in question possess certain symmetries, e.g. translation-invariance on the lattice, and have spatially local interactions. Then when we perform RG flow on the theory, the log-probability (i.e. the negative of the *action* of the theory) will change and possibly gain new terms; the RG flow has the effect of zooming out and coarsening our description of the probability distribution. That is, the RG flow only preserves data from the probability distribution that allows us to compute expectation values of functions which are smooth on large distance scales. What is perhaps a priori surprising is that, in many circumstances, the RG-flowed theory can be well-described by a *finite* number of parameters, which manifest as couplings in the log-probability. This is to say that we can truncate the log-probability to a finite number of terms. Interestingly, there are many distinct microscopic theories which flow to log-probabilities with the same finite set of terms (albeit with coefficients which depend on the microscopic theory); this class of truncated log-probabilities defines a *universality class*. They are akin to the stable distributions in classical probability theory, i.e. the Gaussian distributions, Cauchy distributions, etc. We emphasize that in field theory, the effective truncation comes into play at *intermediate distance scales*; this means that we do not have to RG flow our UV theory until it hits a fixed point. When we have a UV theory that has been RG flowed so that it approximately sits in a universality class, we call the resulting probability distribution an *IR theory*, coming from the terminology of IR light being long-wavelength (relative to visible light).

The universality classes themselves, each specified by a finite number of parameters, are called *Effective Field Theories* or *EFTs* [74,75]. They provide an accurate description of the long-distance behavior of many physically important microscopic theories. In fact, the Standard Model of particle physics is an example of an EFT, which can be understood through the following perspective. Suppose we know the underlying symmetries and types of particles observed at colliders. Then we imagine that the particles we observe are themselves built out of much smaller, more microscopic degrees of freedom which are inaccessible to us. As such, we stipulate that the particles we do observe should be described by a universality class corresponding to the long-distance behavior of the inaccessible, microscopic degrees of freedom. Then we can write down the EFT consistent with the data we can observe; the EFT has a finite number of parameters which can be determined by looking at the detailed experimental data. Remarkably, this procedure yields the Standard Model as we know it.

4.7.2 A heuristic example

To gain some intuition for how EFT can work, let us go back to our favorite example of scalar field theory, here in d spatial dimensions. We will make some heuristic arguments which work equally well in a lattice or continuum formulation; we opt for the continuum formulation since it makes the arguments slightly easier to phrase. We accordingly consider the negative of the log-probability (i.e. the action)

$$\int_{\mathbb{R}^d} dy \left(\frac{1}{2} \nabla \phi \cdot \nabla \phi + \frac{1}{2} m^2 \phi^2 + \frac{\lambda}{4!} \phi^4 \right). \quad (4.36)$$

This theory possesses symmetry under rotations and translations, as well as under $\phi \rightarrow -\phi$. Let us suppose that our RG flow (approximately) respects these symmetries, and as such only generates terms which are (approximately) consistent with said symmetries. If we think of y as corresponding to a position coordinate which has units of length, then dy has units of length to the d th power. We denote this fact by $[dy] = [\ell]^d$, in evident notation. Since the log-probability, e.g. (4.36), is unitless, we can work out the dimensions of ϕ . Since $[\nabla] = [\ell]^{-1}$, it follows that $[\phi] = [\ell]^{-\frac{d-2}{2}}$ in order for $\int_{\mathbb{R}^d} dy \frac{1}{2} \nabla \phi \cdot \nabla \phi$ to be dimensionless. Similarly, we determine that $[m] = [\ell]^{-1}$ and $[\lambda] = [\ell]^{d-4}$. Let us define the dimensionless parameters $\hat{m} := \ell m$ and $\hat{\lambda} := \ell^{4-d} \lambda$, the dimensionless coordinate $\hat{y} = y/\ell$, and the dimensionless field $\hat{\phi} := \ell^{\frac{d-2}{2}} \phi$. In these definitions, ℓ can be regarded as a convenient length scale at which we perform measurements. Then we can rewrite (4.36) as

$$\int_{\mathbb{R}^d} d\hat{y} \left(\frac{1}{2} \widehat{\nabla} \hat{\phi} \cdot \widehat{\nabla} \hat{\phi} + \frac{1}{2} \ell^2 \hat{m}^2 \hat{\phi}^2 + \ell^{4-d} \frac{\hat{\lambda}}{4!} \hat{\phi}^4 \right). \quad (4.37)$$

Before proceeding with our arguments, we note that we have chosen in (4.36) and thus in (4.37) not to put a (dimensionful) coupling in front of the kinetic term $\nabla \phi \cdot \nabla \phi$. This is akin to our previous discussion surrounding (4.12) where we discussed removing factors of $a_{NN'}$ in the kinetic term in the lattice action by rescaling ϕ . In short, we omit the coupling in front of the kinetic term in order to canonically normalize ϕ so that its fluctuations (which are dictated by the kinetic

term) have unit size independent of the length scale at which we measure ϕ ; this is tantamount to a definition of the scalar field ϕ , and in particular calibrates what it means to measure it.

Returning our attention to (4.37), we use the following heuristic. Note that in the equation we have not accounted for any RG flow which would smooth out to a length scale ℓ . However, we can use the ℓ -scaling of the couplings \hat{m} and $\hat{\lambda}$ as a heuristic proxy for how important the corresponding terms in the log-density might be to measurements at length scale ℓ . This heuristic is not entirely reliable, but it often indicates the right answer. Note that if ℓ is a large length, then the $\frac{1}{2} \ell^2 \hat{m}^2 \hat{\phi}^2$ term is large. This indicates that this quadratic term in $\hat{\phi}$ is important at long distances. On the other hand, the importance of the quartic term $\ell^{4-d} \hat{\lambda} \hat{\phi}^4$ depends on the dimension; for $d < 4$ the term is large for large ℓ , for $d = 4$ the term does not scale with ℓ , and for $d > 4$ the term shrinks with increasing ℓ .

More broadly, suppose we write down the most general version of (4.37) that is a sum of polynomials of $\hat{\nabla}\hat{\phi}$'s and $\hat{\phi}$'s, and which satisfies translational and rotational symmetry as well as a symmetry under $\phi \rightarrow -\phi$. Then we would have

$$\int_{\mathbb{R}^d} dy \left(\frac{1}{2} \hat{\nabla}\hat{\phi} \cdot \hat{\nabla}\hat{\phi} + \frac{1}{2} \ell^2 \hat{m}^2 \hat{\phi}^2 + \ell^{4-d} \frac{\hat{\lambda}}{4!} \hat{\phi}^4 + \ell^{6-2d} \hat{\alpha} \phi^6 + \ell^{2-d} \hat{\beta} \hat{\phi}^2 (\hat{\nabla}\hat{\phi} \cdot \hat{\nabla}\hat{\phi}) + \dots \right). \quad (4.38)$$

For $d > 4$, only the first two terms survive for large ℓ ; all of the other terms go to zero. As such, one expects that EFT in $d > 4$ dimensions for the scalar with our chosen symmetries is characterized by a single parameter, namely \hat{m} . For $d = 4$, the only surviving terms for large ℓ are the first three; as such one expects that EFT in $d = 4$ is characterized by the two parameters \hat{m} and $\hat{\lambda}$. In a similar vein, we expect that EFT in $d = 3$ is characterized by the parameters \hat{m} , $\hat{\lambda}$, and $\hat{\alpha}$.

One can check, however, that for $d = 2$, infinitely many parameters are important when ℓ is large; in particular, we can have terms like $\int_{\mathbb{R}^2} dy \ell^2 \hat{\gamma} \hat{\phi}^{2n}$ and $\int_{\mathbb{R}^2} dy \hat{\delta} \hat{\phi}^{2n} (\hat{\nabla}\hat{\phi} \cdot \hat{\nabla}\hat{\phi})$. While there are infinitely many terms, they have a somewhat constrained form. For $d = 1$ we likewise have infinitely many terms which are important for large ℓ , but the types of terms are more varied.

We see from the above heuristics that EFT is most constraining in higher dimensions; in $d \geq 3$ there are only finitely many terms which survive for large ℓ . The EFT paradigm remains somewhat useful even in $d = 2$ where the terms surviving at large ℓ are usefully constrained.

The above analysis can be generalized to incorporate multiple scalar fields or other kinds of fields beyond scalars (e.g. fermions, vector fields, tensor fields, etc.). Moreover, the heuristic analysis can be refined to more carefully account for the effects of RG. One subtlety is that there are circumstances in which RG produces terms in the log-probability which are non-polynomial in the $\hat{\nabla}\hat{\phi}$'s and $\hat{\phi}$'s; however, in many circumstances such non-polynomial terms can be approximated at long distance scales by polynomial terms. Even recognizing such subtleties, the above gives a flavor for the universality enjoyed by a wide variety of theories at long distances (i.e., in the IR).

4.7.3 Concluding remarks

We emphasize that not all UV theories are RG flowed into universality classes described by a finite number of parameters. This means that not all theories fit into the EFT paradigm.

Let us remark that while our account of the EFT perspective of the Standard Model is the contemporary viewpoint, the Standard Model was developed in a more circuitous manner. It was only near the end, or perhaps slightly after, the Standard Model had come into being that its justification in terms of EFT was post facto formulated (see e.g. [74]). This formulation was and is viewed as a grand synthesis.

Indeed, the EFT paradigm is extremely powerful, and has a wide range of applicability including fundamental particle physics (e.g. [74, 75]), condensed matter systems (e.g. [76, 77]), fluids (e.g. [78, 79]), and even the collective behavior of certain kinds of biological systems (e.g. [80, 81]). From a practical point of view, EFT can be thought of as an organization principle that is highly useful for modeling in the circumstances in which it is valid.

5 Renormalizing diffusion models and multiscale modeling

5.1 Overview

The renormalization group organizes physical models into a hierarchy of theories associated to physics at different length scales. It is often the case that a UV theory flows to an IR theory described by EFT, and thus at long distances is approximately parameterized by a finite number of couplings. As previously discussed, these EFTs can accurately capture coarse-grained statistics; in addition, the simplified equations they provide are often efficient to simulate on a computer. While EFTs can be useful, to capture the full range of relevant physical phenomena one must account for additional terms in the log-probability which are suppressed but non-zero when we probe features of the theory at intermediate-to-large distance scales. The challenge of appropriately adapting coarse-grained simulations to take into account corrections arising from more fine-grained physics is the purview of *multiscale modeling* [82].

The subject of multiscale modeling is rather capacious. It includes as a special case EFTs, which as previously discussed are essential to our understanding of particle physics, statistical field theories, and condensed matter systems [7]. Similar phenomena arise in fluid dynamics, where e.g. correction terms to simplified models are needed to understand fusion physics and stellar plasmas. For example, there are important corrections to the ideal MHD equations describing low-frequency, large-scale plasma dynamics, which arise from the Maxwell-Vlasov equations – a higher-dimensional PDE that is comparatively more difficult to simulate [83, 84]. In another example, the large-scale structure of the cosmic web can be understood through the Zeldovich approximation or through Lagrangian Perturbation Theory [85, 86], which in principle is an uncontrolled approximation to cosmological models involving N-body simulations [87, 88]. It is natural to try to learn better coarse-grained approximations from comparatively expensive fine-grained simulations using

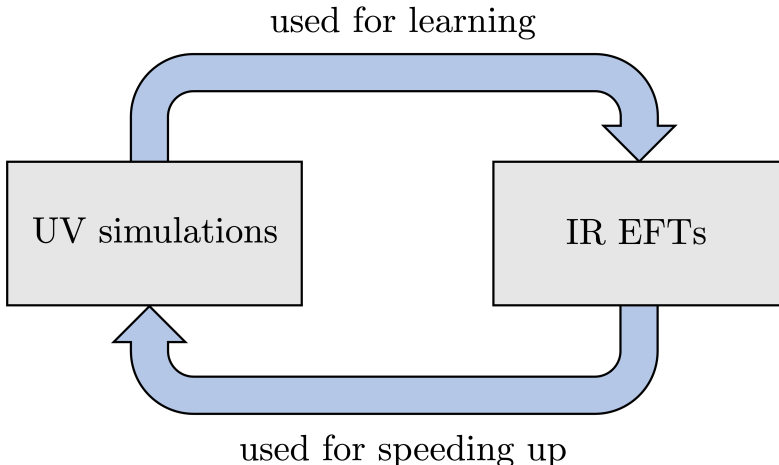


Figure 5: A diagram outlining how our algorithms leverage UV simulations to learn IR EFTs and in turn leverage IR EFTs to speed up UV simulations.

machine-learning (ML) techniques (see e.g. [89,90]).

While coarse-grained and fine-grained physics are rigorously connected through the renormalization group, approaches to improve coarse-grained models using machine learning can lack interpretability and be difficult to evaluate due to a lack of physical interpretation of the model parameters [91]. In our view, it is natural to look for principles for designing ML-based multiscale models for which model parameters are interpretable through the unifying formalism of EFT and the renormalization group, by analogy to existing work on physics-constrained ML models [92,93] which operate on a single scale. In other words, we should be able to rigorously *learn IR EFTs* from UV microscopic simulations, and dually use learned EFTs to *speed up UV simulations*, as schematized in Figure 5.

While we view the development of a general rigorous methodology for multiscale ML-physics simulations as an exciting research problem, in this Section we focus on a concrete setting where the renormalization group formalism is fully developed, and has been studied with ML-based techniques: the problem of *sampling from a statistical field theory*. Specifically, we will turn to the problem of building improved Monte Carlo samplers for densities $\mu_{N,\{\lambda_{i,N}\}}[\phi_N]$ like (4.5) or (4.18) (where $\{\lambda_{i,N}\}$ are the parameters in the action) associated to statistical field theories. In the setting of statistical field theory the log-density is known, and so it is possible to draw samples using the methods reviewed in Section 2. The basic challenge is to design the transition probabilities $Q(\phi|\phi')$ such that the Markov chain mixes well. This becomes very challenging as one makes the lattice parameter $N = L/\epsilon$ large due to the N^d -dependence of the dimension of $\mu_{N,\{\lambda_{i,N}\}}[\phi_N]$.

Various measures of mixing times can have a power-law dependence on N , especially when parameters $\{\lambda_{i,N}\}$ are tuned to match an RG fixed point; this dependence goes by the name of *critical slowing down* [9]. Moreover, in certain kinds of physically relevant models, there are long-distance modes governed by topological quantities like the *topological charge* which contribute in

complex ways to physical processes such as instanton condensation [94,95]; these topological charges can be essential for understanding and characterizing the critical slowing down. Another way to characterize mixing times is to compute the *integrated autocorrelation time* of the sampler, which can be estimated via time-averages of rejection rates. Various sophisticated methods, many of them based on RG, have been used by the lattice field theory community to design samplers which improve on the aforementioned measures [96–100]. Some of the methods are based around bridge sampling and the renormalization group, wherein one introduces several parallel Markov chains for models sampling from $\mu_{N,\{\lambda_i,N\}}[\phi_N]$, such that the bridge used in the parallel tempering algorithm is an approximation to the RG flow of a field theory [96,97].

Alternatively, instead of trying to use physically-inspired ansätze for the proposal distribution $Q_\theta(\phi|\phi')$, one can try to *learn* the proposal distribution using a variational family $Q(\phi|\phi')$. We turn to this perspective below.

5.2 Variationally optimizing the proposal distribution: normalizing flows

In the rest of this Section, we will use the notation ϕ and ϕ_N interchangeably when the context is clear. We will similarly use $d\phi$ and $\mathcal{D}\phi_N$ interchangeably. Let us fix, for the remainder of this subsection, the dimension of the lattice to be N_0 . Suppose we desire to learn the distribution

$$\mu_{N_0,\{\lambda_i,N_0\}}(\phi) := \frac{1}{Z_{\lambda_{N_0}}} e^{-S_{N_0,\{\lambda_i,N_0\}}(\phi)} \mathcal{D}\phi_{N_0} =: p(\phi) d\phi. \quad (5.1)$$

Since the minus log-density $-\log p(\phi)$ is given by the action $S_{N_0,\{\lambda_i,N_0\}}$ of the field theory and thus explicitly known up to a normalizing constant $\log Z_{N_0,\{\lambda_i,N_0\}}$, one can use Metropolis-Hastings to produce samples from this distribution. Given a ψ -independent family of proposal distributions $Q_\theta(\phi|\psi) = Q_\theta(\phi)$, it is natural to try to minimize one of the KL divergences

$$\text{KL}(Q_\theta(\phi)|p(\phi)) \quad \text{or} \quad \text{KL}(p(\phi)|Q_\theta(\phi)) \quad (5.2)$$

over the parameters θ . Optimizing either of these quantities has different trade-offs for Monte Carlo sampling, as optimizing the former tends to find distributions which underestimate the support of the target distribution $p(\phi)$ while optimizing the latter tends to find distributions which overestimate the support of the target distribution [43].

Let us focus for the time being on the problem of optimizing $\text{KL}(Q_\theta|p)$. In this case, a straightforward way to define a family of proposal distributions $Q_\theta(\phi)$ is to (i) fix an an easy-to-sample distribution $p_\infty(\phi') d\phi'$ such as a Gaussian, (ii) define a family of diffeomorphisms $\phi \mapsto \phi' = f_\theta(\phi)$, and (iii) set

$$Q_\theta(\phi) d\phi := (f_\theta^{-1})_*(p_\infty(\phi') d\phi') = p_\infty(f_\theta(\phi)) \left| \det \frac{\partial f_\theta(\phi)}{\partial \phi} \right| d\phi. \quad (5.3)$$

Then so long as the Jacobian factor in the above formula takes a simple analytical form for the parameterization $\theta \mapsto f_\theta$, one can use automatic differentiation and stochastic gradient descent

to minimize the loss

$$L(\theta) = \text{KL}(Q_\theta|p) - \log Z_{N_0, \{\lambda_{i, N_0}\}} = \int d\phi Q_\theta(\phi) \left(\log Q_\theta(\phi) + S_{N_0, \{\lambda_{i, N_0}\}}(\phi) \right). \quad (5.4)$$

From (2.6), this loss admits a straightforward Monte Carlo gradient estimator, since

$$\nabla_\theta L(\theta) = \mathbb{E}_{\phi \sim Q_\theta(\phi)} \left[(\nabla_\theta \log Q_\theta(\phi)) \left(\log Q_\theta(\phi) + S_{N_0, \{\lambda_{i, N_0}\}}(\phi) + 1 \right) \right] \quad (5.5)$$

where as usual we use the fact that $Z_{N_0, \{\lambda_{i, N_0}\}}$ and $S_{N_0, \{\lambda_{i, N_0}\}}$ are θ -independent. We can readily estimate the gradients arising on the right-hand side by sampling from $Q_\theta = (f_\theta^{-1})_*(p_\infty(\phi') d\phi')$: we plug the samples into the analytically computed quantity in the brackets on the right-hand side of (5.5). The machine learning community has produced a wealth of function approximators, such as the *real NVP* flow [101], which have the feature that the Jacobian factor in (5.3) as well as f_θ^{-1} are both efficiently computable.

The approach described above has been explored extensively in a series of papers [2, 31, 102]. In these papers it is shown that on small lattice sizes ($N < 20$) the aforementioned sampler improves upon HMC with respect to various measures of critical slowing down. However, unlike in HMC, there is a pre-training cost to any such flow-based sampler due to the variational inference needed to learn θ , and a systematic quantification of how much pre-training is needed to improve performance is not yet available for reasonable lattice sizes.

Any reasonable parameterization of the family of diffeomorphisms f_θ is built up by iterative compositions, e.g.

$$f_\theta = f_{\theta_T}^T \circ f_{\theta_{T-1}}^{T-1} \circ \dots \circ f_{\theta_1}^1. \quad (5.6)$$

Thus, writing $p_T^\theta = p_\infty$, such a model in principle defines a sequence of distributions p_t^θ for $t = T, \dots, 0$, interpolating between $p_\infty = p_T$ and $p = p_0$, via

$$p_t^\theta = (f_\theta^{t+1})_*^{-1} p_{t+1}^\theta. \quad (5.7)$$

One can take the continuous-time limit and instead use a continuous family of distributions p_t^θ as the pushforwards $(f_\theta^{t,T})_*^{-1} p_\infty$, where $\phi \mapsto f_\theta^{t,T}(\phi)$ is the flow obtained by solving an ODE

$$\frac{d\phi_\tau}{d\tau} = b_\theta(\phi_\tau, \tau) \quad (5.8)$$

from $\tau = t$ to $\tau = T$. In such a framework, one can get estimates of $\log p_t^\theta$ and $\nabla_\theta \log p_t^\theta$ via Neural ODE methods [61]. This class of models has improved the scaling performance with respect N for sampling in ϕ^4 theory [4].

5.3 Learning the renormalization group

One challenge with the class of models discussed above [2, 4, 31] is that the flow p_t^θ learned by the models has no physical meaning, and thus it can be challenging to debug or tune the models using

physical intuition about the field theories being studied. In this subsection, we propose to modify the objective functions such that *the flow learned by the models will be the RG flow of the field theory*. We will describe a more flexible perspective on this class of models in Section 5.6; for now, we will focus on the use of a single RG scheme, and make the smallest possible modification to the objective (5.4).

Recall from equations (4.32) and (4.34) that the Carosso RG scheme is governed by the SDE

$$\partial_t \phi_t(\mathbf{n}) = \Delta \phi(\mathbf{n}) + \eta_t(\mathbf{n}), \quad \phi_0(\mathbf{n}) = \phi(\mathbf{n}), \quad (5.9)$$

where t is a fictitious time parameter associated with scale, Δ is the discrete Laplacian, and η_t is mean-zero Gaussian noise satisfying

$$\mathbb{E}[\eta_t(\mathbf{n})] = 0, \quad \mathbb{E}[\eta_t(\mathbf{n})\eta_s(\mathbf{m})] := \left(\frac{N}{L}\right)^d \delta(t-s) \delta_{\mathbf{n},\mathbf{m}}. \quad (5.10)$$

The probability flow p_t of $p(\phi) =: p_0(\phi)$ under (5.9) can be computed explicitly, as we now explain. Recall the formula (4.26) for the discrete Fourier transform

$$\tilde{\phi}(\mathbf{p}) := \frac{1}{N^d} \sum_{\mathbf{n} \in \mathbb{Z}_N^d} e^{-i\frac{2\pi}{N}\mathbf{p}\cdot\mathbf{n}} \phi(\mathbf{n}) \quad (5.11)$$

which can be computed in time $O((N \log N)^d)$. Then taking the Fourier transform of (5.9) we obtain the equation

$$\partial_t \tilde{\phi}_t(\mathbf{p}) = -|\hat{\mathbf{p}}|^2 \tilde{\phi}_t(\mathbf{p}) + \tilde{\eta}_t(\mathbf{p}), \quad \tilde{\phi}_0(\mathbf{p}) = \tilde{\phi}(\mathbf{p}), \quad (5.12)$$

where we recall the definition (4.27)

$$\hat{p}_i := \frac{2N}{L} \sin\left(\frac{2\pi}{N} \frac{p_i}{2}\right), \quad i = 1, \dots, d, \quad (5.13)$$

so that $|\hat{\mathbf{p}}|^2 = \sum_{i=1}^d \hat{p}_i^2$, and the Gaussian noise satisfies

$$\mathbb{E}[\tilde{\eta}_t(\mathbf{p})] = 0, \quad \mathbb{E}[\tilde{\eta}_t(\mathbf{p})\tilde{\eta}_s(\mathbf{k})] := \frac{1}{L^d} \delta(t-s) \delta_{\mathbf{p},-\mathbf{k}}. \quad (5.14)$$

Notice that (5.12) naturally contains $|\hat{\mathbf{p}}|^2$ instead of $|\mathbf{p}|^2$, and so as such our RG scheme most naturally suppresses modes with large $|\hat{\mathbf{p}}|$. We further observe that (5.12) is diagonal in $\tilde{\phi}(\mathbf{p})$. It will be convenient to use the notation $\tilde{\eta}_t(\mathbf{p}) dt := \Omega dB_t(\mathbf{p})$ where B_t is standard Gaussian noise and $\Omega = \frac{1}{L^d}$. Recall that the stationary distribution of the 1-dimensional Ornstein-Uhlenbeck process

$$dX_t = -\kappa X_t + \sigma dB_t \quad (5.15)$$

is given by

$$X_\infty \sim \mathcal{N}(0, \sigma^2/2\kappa). \quad (5.16)$$

As such, we see that we can sample $\phi(\mathbf{n})$ from $p_\infty = \lim_{t \rightarrow \infty} p_t$ by sampling

$$\operatorname{Re}\{\tilde{\phi}(\mathbf{p})\}, \operatorname{Im}\{\tilde{\phi}(\mathbf{p})\} \sim \mathcal{N}\left(0, \frac{\Omega}{4|\hat{\mathbf{p}}|^2}\right), \quad \mathbf{p} \neq \mathbf{0} \text{ or } (N/2, N/2, \dots, N/2) \quad (5.17)$$

$$\tilde{\phi}(\mathbf{0}), \tilde{\phi}((N/2, N/2, \dots, N/2)) \sim \mathcal{N}\left(0, \frac{\Omega}{2|\hat{\mathbf{p}}|^2}\right), \quad (5.18)$$

as well as using $\phi(-\mathbf{p}) = \phi(\mathbf{p})^*$, and then further applying the inverse of the transformation (5.11). Similarly, after performing a discrete Fourier transform, the transition kernel for (5.12) is also explicitly computable:

$$p_{t,0}(\tilde{\phi}_t|\tilde{\phi}_0) \propto \prod_{\mathbf{p} \in \mathbb{Z}_N^d} \exp\left(-\frac{1}{2} \frac{2|\hat{\mathbf{p}}|^2}{\Omega(1 - e^{-2|\hat{\mathbf{p}}|^2 t})} \left|\tilde{\phi}_t(\mathbf{p}) - \tilde{\phi}_0(\mathbf{p}) e^{-|\hat{\mathbf{p}}|^2 t}\right|^2\right). \quad (5.19)$$

In particular, given samples from p_0 , one can efficiently sample from p_t , and one has an efficient analytic formula for $\nabla_{\phi'} p_{t,0}(\phi'|\phi)$.

In some instances, it is useful to modify the Carosso flow by deforming with a mass term. It is natural to let $\Omega = \frac{2}{L^d}$ and send $|\hat{\mathbf{p}}|^2 \rightarrow |\hat{\mathbf{p}}|^2 + M^2$ so that the kernel above becomes

$$p_{t,0}(\tilde{\phi}_t|\tilde{\phi}_0) \propto \prod_{\mathbf{p} \in \mathbb{Z}_N^d} \exp\left(-\frac{L^d}{2} \frac{|\hat{\mathbf{p}}|^2 + M^2}{(1 - e^{-2(|\hat{\mathbf{p}}|^2 + M^2)t})} \left|\tilde{\phi}_t(\mathbf{p}) - \tilde{\phi}_0(\mathbf{p}) e^{-(|\hat{\mathbf{p}}|^2 + M^2)t}\right|^2\right). \quad (5.20)$$

Note that as $t \rightarrow \infty$, we obtain the probability distribution over $\tilde{\phi}_t$ for a free scalar field with bare mass M . An important feature of (5.20) is that when $\mathbf{p} = \mathbf{0}$, the argument of the exponential is non-zero; this is not true for (5.19). We find in numerical experiments that (5.20) is more numerically stable than (5.19), and so we use (5.20) in our numerical experiments in Section 7.

We now turn our attention to learning distributions p_t^θ which approximate p_t . Accordingly, we should minimize a t -integral of divergences between p_t^θ and p_t . The functional

$$\int_0^T dt \lambda(t) \operatorname{KL}(p_t^\theta|p_t) \quad (5.21)$$

does not work well for this task, because the gradient estimator (5.5) for the term $\operatorname{KL}(p_t^\theta|p_t)$ would require access to $\log p_t$, while we only have sample access to p_t . Instead, we consider the objective

$$\operatorname{KL}(p_0^\theta|p_0) + \int_0^T dt \lambda(t) \operatorname{KL}(p_t|p_t^\theta). \quad (5.22)$$

We can then parameterize p_t^θ for $t = 0, \dots, T$ using a flow-based model (5.7) or a Neural ODE (5.8). This gives us access to the log-densities p_t^θ and their θ -gradients, which is necessary since we need to compute $\nabla_\theta \operatorname{KL}(p_t|p_t^\theta) = \mathbb{E}_{\phi_t \sim p_t}[\nabla_\theta \log p_t^\theta]$. We can then run an algorithm which improves θ while drawing samples from p_0 . For example, let us parameterize the flow as in (5.8), and let

$\theta = (\theta^1, \dots, \theta^M)$. Then we can run the following loop:

initialize $\theta = (\theta^1, \dots, \theta^M)$, ϕ_0 , $i = 0$
if $i < i_{\max}$ **do**
 sample $\phi'_T \leftarrow p_T^\theta = p_\infty$
 set ϕ'_0 by solving (5.8) from $t = T$ to $t = 0$ with initial condition ϕ'_T
 set $\theta_0 \leftarrow \epsilon \left\{ \nabla_\theta \log Q_\theta(\phi'_0) \left(\log Q_\theta(\phi'_0) - S_{N_0, \{\lambda_{i, N_0}\}}(\phi'_0) + 1 \right) \right\}$
 optional: draw ϕ''_0 from an auxiliary distribution $Q'(\phi'' | \phi')$, e.g. simulate several steps of
 Hamiltonian Monte Carlo for p_0 with initial position ϕ'_0 , and then set $\phi''_0 \leftarrow \phi''_0$
 with probability $\alpha = \exp(-S_{N_0, \{\lambda_{i, N_0}\}}(\phi'_0) + S_{N_0, \{\lambda_{i, N_0}\}}(\phi_0)) = p(\phi'_0)/p(\phi_0)$, set $\phi_0 \leftarrow \phi'_0$, replace
i $\rightarrow i + 1$
 for $j = 1, \dots, \tau$
 sample $t_j \leftarrow [0, T]$
 sample $\phi_j \leftarrow p_{t_j}$ using (5.19)
 compute $\nabla_\theta \log p_{t_j}^\theta(\phi_j)$ by calling an ODE solver, namely:
 (i) solve $dx/dt = b_\theta(x, t)$ from $t = t_j$ to $t = T$ with initial condition $x_t = \phi_t$
 (ii) set $x_T \leftarrow x(T)$
 solve the following system from $t = T$ to $t = t_j$ for $k = 1, \dots, M$:

$$\frac{dx}{dt} = b_\theta(x, t), \quad \frac{d(\delta_{\theta^k} x)}{dt} = \frac{\partial}{\partial \theta^k} b_\theta(x, t), \quad \frac{d(\delta \theta^k)}{dt} = \nabla_x \cdot \left(\frac{\partial}{\partial \theta^k} b_\theta(x, t) \right) + \sum_\ell \sum_m b^\ell(x, t) \delta_{\theta^m} x,$$

 with initial conditions $\delta_{\theta^k} x(T) = 0$, $\delta \theta^k(T) = 0$ for all k where $\delta \theta^k(t) := \frac{\partial}{\partial \theta^k} \log p_t^\theta(x(t))$
 set $\theta_j \leftarrow (\delta \theta^1(t_j), \dots, \delta \theta^M(t_j))$
 end for
 set $\theta \leftarrow \theta - \sum_{j=0}^{\tau} \theta_j$
end if
return θ

Having completed the above algorithm, we can now sample ϕ 's from p_T^θ using our outputted θ . In the course of running the algorithm, we have produced a number of samples of ϕ , which we can also leverage if we so desire (with the understanding that convergence of expectation values with respect to ϕ 's sampled earlier on in the algorithm may have a higher variance).

5.4 Variational inference for sampling and diffusion models

5.4.1 Description of the algorithm

A major downside of the objective (5.22) is that it requires the computation of $\nabla_\theta \log p_t^\theta$ for optimization, which can be expensive as it involves repeated solution of a differential equation. Instead, taking inspiration from Section 3, we notice that we are trying to learn the inverse process to the dif-

fusion process (5.9) which implements the Carosso RG scheme. Thus, instead of minimizing (5.22), we can minimize a weighted sum of score function objectives like (3.27):

$$\begin{aligned} L(\theta) &= \int_0^T dt \lambda(t) \mathbb{E}_{\phi_t \sim p_t} \left[|\nabla \log p_t(\phi_t) - \nabla \log p_t^\theta(\phi_t)| \right] \\ &= \int_0^T dt \lambda(t) \mathbb{E}_{\phi_0 \sim p_0} \mathbb{E}_{\phi_t \sim p_{t,0}(\phi_t|\phi_0)} \left[|\nabla_{\phi_t} \log p_{t,0}(\phi_t|\phi_0) - \nabla_{\phi_t} \log p_t^\theta(\phi_t)| \right] + C \end{aligned} \quad (5.23)$$

for a constant C [60].

The rewriting of the objective in the second line is helpful because $\nabla_{\phi_t} \log p_{t,0}(\phi_t|\phi_0)$ is explicit due to (5.19), and the optimization of such an objective can be parallelized across different values of t . Moreover, with the above objective, one only needs access to the score functions $s_t^\theta(\phi_t) = \nabla_{\phi_t} \log p_t^\theta(\phi_t)$, which can now be parameterized using a neural net architecture like a U-Net [58] that have been developed for highly-successful applications in the generative modeling of images. Writing (5.9) in the notation

$$d\tilde{\phi}_t = -\Delta\phi dt + \sigma dB_t \quad (5.24)$$

where for us $\sigma = 1$, the corresponding time-reversed process is

$$d\tilde{\phi}_t = (-\Delta\phi + \sigma^2 s_t[\phi_t]) dt + \sigma dB_t \quad (5.25)$$

where $s_t = \nabla_\phi \log p_t$ is the score function. In this setting, the training algorithm is as follows:

initialize θ, ϕ_0

if $i < i_{\max}$ **do**

sample $\phi_T \leftarrow p_T^\theta = p_\infty$

set ϕ'_0 by solving (5.25) from $t = T$ to $t = 0$ with initial condition ϕ_T and substituting $s_t \rightarrow s_t^\theta$

optional: draw ϕ''_0 from an auxiliary distribution $Q'(\phi''|\phi')$, e.g. simulate several steps of

Hamiltonian Monte Carlo for p_0 with initial position ϕ'_0 , and then set $\phi'_0 \leftarrow \phi''_0$

with probability $\alpha = \exp(-S_{N_0, \{\lambda_{i, N_0}\}}(\phi'_0) + S_{N_0, \{\lambda_{i, N_0}\}}(\phi_0)) = p(\phi'_0)/p(\phi_0)$ set $\phi_0 \leftarrow \phi'_0$, replace

$i \rightarrow i + 1$

for $j = 1, \dots, \tau$

sample $t_j \leftarrow [0, T]$

sample ϕ_j from p_{t_j} using (5.19)

set $\theta_j \leftarrow \lambda(t_j) (\nabla_{\phi'} p_{t_j}(\phi'|\phi_0) - \nabla_{\phi'} s_{t_j}^\theta(\phi', t_j))|_{\phi'=\phi_j}$

end for

set $\theta \leftarrow \sum_{j=1}^\tau \theta_j$

end if

return θ

Below we will discuss some aspects of the above method in further detail.

5.4.2 Comments on the method

Since we have provided a formal connection between a specific diffusion model and a particular RG scheme (i.e. the Carosso scheme), our method produces a physically interpretable class of ML architectures for sampling from field theories. We elaborate more on variations of this method in Section 5.6, including generalizations to RG schemes other than Carosso’s. Before pursuing such generalizations, there are several conceptual points about our method that are already evident.

While we have initially focused on the Carosso scheme due to the simplicity of its SDE formulation, it has the seemingly unusual property that the induced RG flow for long times maps all probability functionals to a fixed functional with a *fixed* cutoff at the lattice scale. This may seem in tension with the usual discussion about (Wilsonian-type) RG schemes wherein IR fixed points of certain field theories do not have to be free (i.e. Gaussian). The way to resolve this apparent tension is to note that in the Carosso scheme, the convergence rate of the Fourier transform of an arbitrary initial distribution to the final distribution p_∞ is *frequency-dependent*, and thus by performing an appropriate *field renormalization* one can extract a modified flow which does not have to converge to a trivial (Gaussian) fixed point but instead can converge to nontrivial RG fixed points. We further discuss field renormalization and how to incorporate it into our modeling framework in Section 5.6.

5.4.3 Towards score-based EFT

In our method, we have identified the score function s_t^θ with the negative ϕ -gradient of the effective action. Here t parameterizes the cutoff scale Λ_t . In physics applications, much is known about the most important terms contributing to the effective action by using the methodologies of EFT, and so this information can be incorporated explicitly into the structure of the neural networks parameterizing s_t^θ . Note that the reverse SDE (5.25) is a non-parameteric version of the inverse of the renormalization group equation and involves s_t^θ . It may be convenient to parameterize the score function by e.g.

$$s_t^\theta(\phi) = \left(\sum_{i=1}^r f_i(\theta^i, t) g_i(\phi) \right) + \tilde{s}_t^{\tilde{\theta}}(\phi) \quad (5.26)$$

where the right-hand side can be understood as follows:

- $\theta = (\theta_1, \dots, \theta_r, \tilde{\theta})$ are the parameters of the neural network.
- The terms $g_i(\phi)$ are the ϕ gradients of the most significant terms in the effective action as dictated by EFT. It is natural to take, for example, $g_1(\phi) = \Delta\phi$, which is proportional to the ϕ gradient of kinetic terms in the effective lattice action. Here Δ is the discrete Laplacian (4.20). It is natural to further take $g_2(\phi) = \phi$, which is the ϕ gradient of the mass term in the effective action; and similarly $g_3(\phi) = \phi^3$ which is proportional to the ϕ gradient of the ϕ^4 term in the effective action. We emphasize that ϕ should be viewed as a N^d -dimensional vector, and similarly $g_1(\phi)$, $g_2(\phi)$, and $g_3(\phi)$ define N^d -dimensional vectors.

- The terms $f_i(\theta^i, t)$ are scalar-valued neural networks. Upon training, these terms will *learn the coefficients* of the terms g_i in the effective action. As such, the estimated functions f_i are physically meaningful, and give qualitatively novel estimators of fundamental physical quantities, e.g. the bare couplings at scale Λ and the wave function renormalization.
- The term $\tilde{s}_t^{\tilde{\theta}}(\phi)$ is a general neural network which captures the remaining terms needed to describe the (gradient of the) effective action. To decouple the training of $\tilde{s}_t^{\tilde{\theta}}(\phi)$ from the training of the $f_i(\theta^i, t)$'s, it is natural to require that for all θ and t , the vector field $\tilde{s}_t^{\tilde{\theta}}(\phi)$ is orthogonal to the vector fields $g_i(\phi)$, $i = 1, \dots, r$, a condition that is easy to implement when designing the neural network $\tilde{s}_t^{\tilde{\theta}}$ by adding a fixed linear projection as the last layer.

Renormalizing diffusion models can thus be designed to naturally *estimate the flows of parameters of the effective action*. This provides both a new class of estimators for these important physical quantities, as well as a series of heuristics for further designing the architectures and optimization processes for the neural networks f_i and $\tilde{s}_t^{\tilde{\theta}}$, since much is known about the coefficients of the effective action as well as the behavior of the error term under RG flow of many field theories.

Because the (neural network) quantities f_i have a physical interpretation, plots of their values during training may be used as diagnostics for the training, to see e.g. if the training is being slowed down by an RG critical point. Moreover, since phase transitions and the onset of spontaneous symmetry breaking are meant to be captured by the behavior of critical points of the effective action, there are natural methods which use the estimated quantities f_i to detect phase transitions in lattice field theories, which can otherwise be a difficult problem requiring searches for appropriate order statistics. Finally, our method may be combined with transfer learning methods [4], which have been found to speed up training, to learn a single score function or effective action estimator for field theories with a whole range of UV parameters at once. This opens up the exciting possibility of learning a single representation of the entire RG flow diagram of a field theory, or even to automatically search for phase transitions in field theories in a physically justified way using ML techniques.

We will elaborate on our above score-based schemes for EFT and provide numerical examples in [103].

5.5 More history

We now place the discussion of the previous Subsection into a broader context. The first observation to be made is that the noising process (5.24) as well as the inverse process (5.25) already appear in the machine learning literature for image generative modeling under the name of *Blurring Diffusion Models* [18]. More generally, these models fits into a large class of models [19–22, 26], which we call *multiscale diffusion models*, which all share the feature that they generate images in a multi-step process, where earlier steps generate a *coarse approximation to the image* and then later steps *refine the coarse approximation* to produce the final, high-fidelity image. The first motivation for considering such models for image generation is that images naturally have features across a

hierarchy of scales (in fact, the power frequency spectrum of natural image distributions tends to follow power laws [104]), and thus it seems natural to allocate separate model variables for feature generation at different scales. The second motivation is *computational*: it has been found that repeatedly ‘upsampling’ the image is simply computationally more efficient than trying to generate a high-resolution image from scratch. Leveraging the multiscale structure of image models has been found to consistently improve model performance [19, 22].

In fact, the original paper on diffusion modeling [23] is loosely motivated by ideas about the renormalization group. While there has been much work aimed at connecting specific neural network architectures to special cases of RG transformations [105], as well as works that draw heuristic, non-physically-grounded analogies [106–111], the introduction of multiscale diffusion models was done as a practical matter, without connection to physics. In our method detailed in the previous Subsection, we showed how to rigorously interpret an existing diffusion model architecture (when applied to a problem in lattice field theory) as a pre-existing RG scheme, namely the Carosso scheme. In Section 5.6 below, we expand on this connection, indicating how one can *design* diffusion models to implement different classes of RG schemes, and explaining how to compare the results from different multiscale diffusion models when applied to field-theoretic problems by identifying the appropriate rescaling transformations.

The Carosso scheme, which has been our focus so far, is connected to an interesting series of works in the lattice field theory community. Specifically, Carosso [8, 30] introduces his RG scheme as a formal renormalization group counterpart of a widely-used method dubbed “gradient flow”, which descends⁵ from Lüscher’s pioneering [27]. This latter paper pointed out that when trying to sample from lattice $SU(N)$ gauge theory, which is a distribution μ over a compact manifold, there is an implicit characterization of a flow f_t called the *Wilson Flow* such that $(f_t)_*\mu$ limits to a uniform (Haar measure) distribution. In fact, this flow ends up being a lattice analog of the gradient flow of the Yang-Mills functional [28]. As such, it is a PDE for the fields with a gradient flow interpretation which has as its highest order term the Laplacian operator acting on the fields. Evidently the “gradient flow” smooths out the fields and gives rise to an operation akin to the renormalization group; this has led to a significant amount of numerical work in lattice field theory [113], including a widely used method for scale setting [114]. Carosso [8] introduces his scheme as a stochastic analog of the gradient flow, such that certain long-range correlators for the scheme can be computed directly without the stochastic component of the renormalization scheme. In Section 5.6.1, we show that the Carosso scheme [8] can be written in the framework of the Wegner-Morris equation [115–118], and thus put on a common footing with more familiar schemes such as Polchinski renormalization [34].

The methods described in the present paper unify the ideas behind Lüscher’s gradient flow, normalizing flows for field theories, and variational characterizations of the renormalization group [119] into a framework for building physically-interpretable ML models for field-theoretic applications.

⁵The earlier-cited work on normalizing flows for lattice field theory [2, 4, 31] is also loosely motivated by the idea of learning Lüscher’s trivializing flows; explicit statements to this end can be found in [112].

Our work also has synergies with the framework of “Bayesian Renormalization” developed in [120], and may also have applications in the related mathematical subject of *stochastic localization* which is reviewed in e.g. [121, 122].

5.6 Other RG schemes

5.6.1 Wegner-Morris RG flow on the lattice

So far we have focused on Carosso’s RG scheme (see e.g. (5.9) and (5.10)) for its simplicity. However, there are a wide variety of other RG schemes available which have different features and tradeoffs.

First we must comment on what qualifies as an RG scheme. Since RG as a conceptual framework is rather capacious, there are few general rules for which RG schemes are allowed versus disallowed; the praxis is to be inclusive of any scheme which produces results which can be justified on the grounds of quantitative reasoning and conceptual appeals to universality. However, there do exist desirable properties for RG flows to possess, which can be achieved by special subclasses of RG schemes. Let us give an example of a desirable property following Wilson [123] and Polchinski [34].

Suppose we have a distribution $p_0(\phi)$ on which we would like to enact an RG flow. On the lattice, the natural short-distance cutoff scale is $\epsilon = L/N$; then we say that the initial probability density $p_0(\phi)$ captures correlations at distances larger than L/N . Equivalently, we can say that $p_0(\phi)$ captures correlations at momenta smaller than $\sim 1/\epsilon = N/L$. For convenience, let us take $\Lambda_0 = 2\sqrt{d}N/L$. Now let $f(t)$ be a strictly monotonically decreasing function of t such that $f(0) = 1$. We desire an RG flow such that $p_t(\phi)$ captures all correlations at momenta smaller than $\Lambda_0 f(t)$. In other words, as we flow $p_t(\phi)$ for increasing t , it continues to capture correlations at momentum scales smaller than $\Lambda_0 f(t)$, but correlations at larger momenta may not be preserved. This is the sense in which the RG flow preserves only correlations at progressively smaller momentum scales, corresponding to progressively larger distance scales.

More precisely, we want our RG flowed distribution $p_t(\phi)$ to satisfy

$$\mathbb{E}_{\phi \sim p_t(\phi)}[\tilde{\phi}(\mathbf{p}_1)\tilde{\phi}(\mathbf{p}_2)\cdots\tilde{\phi}(\mathbf{p}_r)] \approx \mathbb{E}_{\phi \sim p_0(\phi)}[\tilde{\phi}(\mathbf{p}_1)\tilde{\phi}(\mathbf{p}_2)\cdots\tilde{\phi}(\mathbf{p}_r)], \quad \text{for all } |\hat{\mathbf{p}}_i| \leq \Lambda_0 f(t) \quad (5.27)$$

for some strictly monotonically decreasing $f(t)$ with $f(0) = 1$. Above, the number r of $\tilde{\phi}$ ’s in the correlator is arbitrary. The property (5.27) means that $p_t(\phi)$ has the same long-distance correlators (dictated by small momentum) as $p_0(\phi)$. As t grows, the momentum scale below which correlators of p_0 and p_t agree becomes smaller and smaller; equivalently, the distance scale above which p_0 and p_t agree becomes larger and larger. We would also like for $p_t(\phi)$ to induce a flow on samples $\tilde{\phi}(\mathbf{p})$ which mixes modes in momentum space, as discussed in detail in Section 4.6. There is a very nice class of RG schemes which satisfy the property (5.27). One of the more famous examples is Polchinski’s scheme [34], which can be viewed as a special case of the Wegner-Morris flow equation [115–118] (which is discussed in more detail in Appendix C). Indeed, the Carosso scheme [8] can also be viewed as a special case of the Wegner-Morris flow equation.

Let us briefly explain the lattice version of the Wegner-Morris flow equation here. We will specialize to a particular class of Wegner-Morris flows reviewed in e.g. [14, 119]. To define the flow, we require two basic objects, which are prescribed functions of t :

1. A *cutoff function* $\tilde{B}_t(|\hat{\mathbf{p}}|)$. There is a non-increasing function $g(t)$ for $t \geq 0$ with $g(0) = 1$ such that $B_t(|\hat{\mathbf{p}}|)$ goes rapidly to zero for $|\hat{\mathbf{p}}| \geq g(t)$. We further require that $\tilde{B}_t(|\hat{\mathbf{p}}|)$ is $O(1)$ for $|\hat{\mathbf{p}}| \leq \Lambda_0 g(t)$. Note that this $g(t)$ is distinct from the $f(t)$ discussed above; in fact we will later see that $f(t) \leq g(t)$.
2. A *seed probability density* $q_t(\phi)$. This probability density has the property that

$$\mathbb{E}_{\phi \sim q_t(\phi)}[\tilde{\phi}(\mathbf{p}_1) \tilde{\phi}(\mathbf{p}_2) \cdots \tilde{\phi}(\mathbf{p}_r)] \approx 0 \quad \text{for any } |\hat{\mathbf{p}}_i| \geq \Lambda_0 g(t), \quad (5.28)$$

where $g(t)$ is the same function which controls the cutoff function. This means that correlation functions with momenta $|\hat{\mathbf{p}}|$ greater than $\Lambda_0 g(t)$ are suppressed. By convention we write the seed probability density as $q_t(\phi) = \frac{1}{Z_q} e^{-2\hat{S}_t(\phi)}$, where $\hat{S}_t(\phi)$ is called the *seed action*.

With the above ingredients, the lattice Wegner-Morris flow equation is

$$\partial_t p_t(\phi) = \frac{1}{2} \sum_{\mathbf{n}, \mathbf{m} \in \mathbb{Z}_N^d} B_t(|\mathbf{n} - \mathbf{m}|_{\mathbb{Z}_N^d}) \left(\frac{\partial^2 p_t(\phi)}{\partial \phi(\mathbf{n}) \partial \phi(\mathbf{m})} + 2 \frac{\partial}{\partial \phi(\mathbf{n})} \left(\frac{\partial \hat{S}_t(\phi)}{\partial \phi(\mathbf{m})} p_t(\phi) \right) \right) \quad (5.29)$$

where we have used

$$B_t(|\mathbf{n} - \mathbf{m}|_{\mathbb{Z}_N^d}) := \sum_{\mathbf{p} \in \mathbb{Z}_N^d} e^{i \frac{2\pi}{N} \mathbf{p} \cdot (\mathbf{n} - \mathbf{m})} \tilde{B}_t(|\mathbf{p}|). \quad (5.30)$$

We notice that (5.29) is a lattice version of a *convection-diffusion equation*, where $\frac{\partial^2 p_t(\phi)}{\partial \phi(\mathbf{n}) \partial \phi(\mathbf{m})}$ is the diffusive term and $2 \frac{\partial}{\partial \phi(\mathbf{n})} \left(\frac{\partial \hat{S}_t(\phi)}{\partial \phi(\mathbf{m})} p_t(\phi) \right)$ is the convective term. We can rewrite (5.29) in momentum space as

$$\partial_t p_t(\tilde{\phi}) = \frac{1}{2} \sum_{\mathbf{p} \in \mathbb{Z}_N^d} \tilde{B}_t(|\mathbf{p}|) \left(\frac{\partial^2 p_t(\tilde{\phi})}{\partial \tilde{\phi}(\mathbf{p}) \partial \tilde{\phi}(-\mathbf{p})} + 2 \frac{\partial}{\partial \tilde{\phi}(\mathbf{p})} \left(\frac{\partial \hat{S}_t(\tilde{\phi})}{\partial \tilde{\phi}(-\mathbf{p})} p_t(\tilde{\phi}) \right) \right) \quad (5.31)$$

It is pleasing that the Wegner-Morris flow is an RG flow scheme taking the form of a convection-diffusion equation for $p_t(\phi)$. This makes intuitive sense: a convection-diffusion equation smooths out a probability density. Interestingly, our requirements on the cutoff function and seed probability density in fact give the convection-diffusion equation special structure, in particular so that (5.27) is satisfied. The connection between RG flows and convection-diffusion equations, as well as their direct interplay with the theory of optimal transport, is analyzed in detail in [119]. A recent study [120] further connects these kinds of convection-diffusion equations to the subject of Bayesian inference.

It may appear that (5.29) is rather different than e.g. the Carosso scheme, which we formulated using an SDE for ϕ_t . However, there is in fact an SDE formulation of (5.29), which as a special case reproduces the Carosso scheme. In particular, consider the SDE

$$\partial_t \phi_t(\mathbf{n}) = - \sum_{\mathbf{m} \in \mathbb{Z}_N^d} B_t(|\mathbf{n} - \mathbf{m}|_{\mathbb{Z}_N^d}) \frac{\partial \widehat{S}_t(\phi)}{\partial \phi(\mathbf{m})} + \eta_t(\mathbf{n}), \quad \phi_0(\mathbf{n}) = \phi(\mathbf{n}), \quad (5.32)$$

with the noise being Gaussian and satisfying

$$\mathbb{E}[\eta_t(\mathbf{n})] = 0, \quad \mathbb{E}[\eta_t(\mathbf{n})\eta_s(\mathbf{m})] = \delta(t-s) B_t(|\mathbf{n} - \mathbf{m}|_{\mathbb{Z}_N^d}). \quad (5.33)$$

Letting $\phi_t[\psi, \eta_t]$ be a solution (5.32) with initial condition ψ and for a fixed sample of the noise η_t , the Wegner-Morris flow equation (5.29) is equivalent to

$$p_t(\phi) = \int \prod_{\mathbf{n} \in \mathbb{Z}_N^d} d\psi(\mathbf{n}) \mathbb{E}_{\eta_t}[\delta(\phi - \phi_t[\psi, \eta_t])] p_0(\psi). \quad (5.34)$$

In other words, to sample from a $p_t(\phi)$ at time t which satisfies the Wegner-Morris equation (5.32), we simply have to sample $\phi \leftarrow p_0(\phi)$ and flow ϕ to time t according to the SDE given by (5.32), (5.33). This idea has been emphasized recently in the literature by Carosso [8] in the context of his scheme.

For completeness, we include the momentum space version of (5.32) and (5.33), namely

$$\partial_t \widetilde{\phi}_t(\mathbf{p}) = -\widetilde{B}_t(|\widehat{\mathbf{p}}|) \frac{\partial \widehat{S}_t(\widetilde{\phi})}{\partial \widetilde{\phi}(-\mathbf{p})} + \widetilde{\eta}_t(\mathbf{p}), \quad \widetilde{\phi}_0(\mathbf{p}) = \widetilde{\phi}(\mathbf{p}), \quad (5.35)$$

where the noise is Gaussian and satisfies

$$\mathbb{E}[\widetilde{\eta}_t(\mathbf{p})] = 0, \quad \mathbb{E}[\widetilde{\eta}_t(\mathbf{p})\widetilde{\eta}_s(\mathbf{k})] = \delta(t-s) \widetilde{B}_t(|\widehat{\mathbf{p}}|) \delta_{\mathbf{p}, -\mathbf{k}}. \quad (5.36)$$

Below we explain how the Wegner-Morris equation (5.32), or equivalently its SDE formulation in (5.32), (5.33), specializes to the Carosso scheme and the Polchinski scheme. Then we make more comments about Wegner-Morris flows more generally, including how to compare between different RG schemes.

5.6.2 Recovering the Carosso scheme

By comparing (5.32), (5.33) with (5.9), (5.10), we see that the Carosso scheme corresponds to

$$B_t(|\mathbf{n} - \mathbf{m}|_{\mathbb{Z}_N^d}) = \left(\frac{N}{L}\right)^d \delta_{\mathbf{n}, \mathbf{m}}, \quad \widehat{S}_t(\phi) = -\frac{(L/N)^d}{2} \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \phi(\mathbf{n}) \Delta \phi(\mathbf{n}). \quad (5.37)$$

In particular, both the cutoff function B_t and the seed action \widehat{S}_t (coming from the seed probability density $q_t(\phi)$) are t -independent in this scheme. This corresponds to the special case of $g(t) = 1$.

In momentum space, the Carosso scheme suppresses the mean value of the $\widetilde{\phi}(\mathbf{p})$'s when $|\widehat{\mathbf{p}}|$ is sufficiently large; in particular from (5.19) we saw that $\mathbb{E}_{\eta_t}[\widetilde{\phi}_t[\widetilde{\phi}_0, \eta_t](\mathbf{p})] = e^{-|\widehat{\mathbf{p}}|^2 t} \widetilde{\phi}_0(\mathbf{p})$. As such, examining (5.19), we approximately have $f(t) \sim \frac{L}{2N\sqrt{d}} \frac{1}{\sqrt{t}}$ for large t .

5.6.3 Recovering the Polchinski scheme

In [34], Polchinski designed a nice RG scheme for a scalar ϕ^4 theory in the continuum. Here we give a latticized version of Polchinski's flow. Polchinski begins by writing out the scalar ϕ^4 theory with a cutoff in momentum space, which we recapitulate on the lattice. Defining $K_t(|\hat{\mathbf{p}}|) := e^{-|\hat{\mathbf{p}}|^2/(e^{-t}\Lambda_0)^2}$, we have

$$p_0(\phi) = \frac{1}{Z} \exp \left(-L^d \sum_{\mathbf{p} \in \mathbb{Z}_N^d} \frac{1}{K_0(|\hat{\mathbf{p}}|)} \frac{1}{2} \tilde{\phi}(\mathbf{p})(|\hat{\mathbf{p}}|^2 + m^2) \tilde{\phi}(-\mathbf{p}) \right. \\ \left. - L^d \sum_{\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3, \mathbf{p}_4 \in \mathbb{Z}_N^d} \frac{\lambda}{4!} \tilde{\phi}(\mathbf{p}_1) \tilde{\phi}(\mathbf{p}_2) \tilde{\phi}(\mathbf{p}_3) \tilde{\phi}(\mathbf{p}_4) \delta_{\mathbf{p}_1 + \mathbf{p}_2 + \mathbf{p}_3 + \mathbf{p}_4, 0} \right), \quad (5.38)$$

where $\Lambda_0 = 2\sqrt{d}N/L$ as usual. Notice the $\frac{1}{K_0(|\hat{\mathbf{p}}|)} = e^{|\hat{\mathbf{p}}|^2/\Lambda_0^2}$ term appearing in front of the $\frac{1}{2} \tilde{\phi}(\mathbf{p})(|\hat{\mathbf{p}}|^2 + m^2) \tilde{\phi}(-\mathbf{p})$ part of the action. At the moment it may appear superfluous; it forces modes $\tilde{\phi}(\mathbf{p})$ with $|\hat{\mathbf{p}}| \gtrsim \Lambda_0$ to be zero with small variance, but this is already accomplished by the lattice discretization since momenta with $|\hat{\mathbf{p}}| > \Lambda_0$ simply do not exist. This being said, the role of the $\frac{1}{K_0(|\hat{\mathbf{p}}|)} = e^{|\hat{\mathbf{p}}|^2/\Lambda_0^2}$ term will become clear shortly.

Considering the Polchinski scheme, in momentum space we have

$$\tilde{B}_t(|\hat{\mathbf{p}}|) = -\frac{1}{L^d} \frac{1}{|\hat{\mathbf{p}}|^2 + m^2} \partial_t K_t(|\hat{\mathbf{p}}|), \quad \tilde{S}_t(\phi) = \frac{L^d}{2} \sum_{\mathbf{p} \in \mathbb{Z}_N^d} \frac{1}{K_t(|\hat{\mathbf{p}}|)} \tilde{\phi}(\mathbf{p})(|\hat{\mathbf{p}}|^2 + m^2) \tilde{\phi}(-\mathbf{p}). \quad (5.39)$$

This scheme is designed so that the RG flow of $p_0(\phi)$ for $\lambda = 0$ is simply

$$p_t(\phi) = \frac{1}{Z_t} \exp \left(-\frac{L^d}{2} \sum_{\mathbf{p} \in \mathbb{Z}_N^d} \frac{1}{K_t(|\hat{\mathbf{p}}|)} \tilde{\phi}(\mathbf{p})(|\hat{\mathbf{p}}|^2 + m^2) \tilde{\phi}(-\mathbf{p}) \right). \quad (5.40)$$

That is, the *free theory* (i.e. having a Gaussian initial probability density $p_0(\phi)$ since $\lambda = 0$) has a simple flow, in which modes with $|\hat{\mathbf{p}}| \gtrsim e^{-t}\Lambda_0$ are exponentially suppressed. In a sense, this nice feature of the Polchinski flow is the reason for the $\frac{1}{K_0(|\hat{\mathbf{p}}|)}$ term in (5.38); this term conspires with (5.39) to produce the simple flow in (5.40) of the free theory.

Even for general λ , the Polchinski flow (with the particular choice of $K_t(|\hat{\mathbf{p}}|)$ given here), the Polchinski scheme suppresses the mean value of $\tilde{\phi}(\mathbf{p})$'s as $\mathbb{E}_{\eta_t}[\tilde{\phi}_t[\tilde{\phi}_0, \eta_t](\mathbf{p})] = \exp\left(-e^{2t} - 1\right) \frac{|\hat{\mathbf{p}}|^2}{\Lambda_0^2} \tilde{\phi}_0(\mathbf{p})$. Accordingly, we have that $f(t) = g(t) = e^{-t}$. Thus the notion of time t in the Polchinski scheme is exponentially different than in the Carosso scheme: the Polchinski scheme suppresses high-momentum modes exponentially faster.

The transition kernel $\mathbb{E}_{\eta_t}[\delta(\phi_t - \phi_t[\phi_0, \eta_t])]$ in the Polchinski setting is given by

$$p_{t,0}^{\text{Polchinski}}(\tilde{\phi}_t|\tilde{\phi}_0) \propto \prod_{\mathbf{p} \in \mathbb{Z}_N^d} \exp \left(-\frac{L^d}{2} \frac{|\hat{\mathbf{p}}|^2 + m^2}{K_t(|\hat{\mathbf{p}}|^2) - \frac{K_t(|\hat{\mathbf{p}}|^2)^2}{K_0(|\hat{\mathbf{p}}|^2)}} \left| \tilde{\phi}_t(\mathbf{p}) - \frac{K_t(|\hat{\mathbf{p}}|^2)}{K_0(|\hat{\mathbf{p}}|^2)} \tilde{\phi}_0(\mathbf{p}) \right|^2 \right) \quad (5.41)$$

This expression holds for any $K_t(|\widehat{\mathbf{p}}|^2)$, i.e. not merely $K_t(|\widehat{\mathbf{p}}|^2) = e^{-|\widehat{\mathbf{p}}|^2/(e^{-t}\Lambda_0)^2}$. For purposes of slowing down the flow, we will later use the kernel

$$K_t(|\widehat{\mathbf{p}}|^2) = e^{-(|\widehat{\mathbf{p}}|^2 + M^2)t} \quad (5.42)$$

where M is a mass parameter satisfying $M/\Lambda_0 \ll 1$. The main advantage of (5.42) is that the exponential of (5.41) does not blow up for $\mathbf{p} = \mathbf{0}$. For our previously-defined kernel this blow-up does occur, which means that the momentum zero mode is preserved by the flow. Using (5.42) is more numerically stable, and so we will opt to use it in Section 7.

5.6.4 More general features

Any Wegner-Morris type schemes, which satisfy (5.27), can be readily compared at low-momentum. In particular, suppose we have two such schemes with different $f(t)$'s, which we will call $f_1(t)$ and $f_2(t)$. Given an initial probability distribution $p_0(\phi)$, suppose its flow under the first scheme is denoted by $p_t^{(1)}(\phi)$, and its flow under the second scheme is denoted by $p_t^{(2)}(\phi)$. Then we have

$$\mathbb{E}_{\phi \sim p_t^{(1)}(\phi)}[\tilde{\phi}(\mathbf{p}_1)\tilde{\phi}(\mathbf{p}_2)\cdots\tilde{\phi}(\mathbf{p}_r)] \approx \mathbb{E}_{\phi \sim p_t^{(2)}(\phi)}[\tilde{\phi}(\mathbf{p}_1)\tilde{\phi}(\mathbf{p}_2)\cdots\tilde{\phi}(\mathbf{p}_r)], \quad \text{for all } |\widehat{\mathbf{p}}_i| \leq \Lambda_0 \min\{f_1(t), f_2(t)\}. \quad (5.43)$$

The minimum on the right-hand side accounts for the fact that the two schemes perform RG at different rates as a function of t . This is all fine, but the more interesting setting for comparison is when both flows reach the same RG fixed point and we compare their properties at *all* momentum scales. To understand this, we need to first address how in the Wegner-Morris flow schemes we can find an RG fixed point in the first place.

Let us begin by discussing the Carosso and Polchinski schemes in particular. In both the Carosso and Polchinski schemes, as t increases, the flow of $p_t(\phi)$ erases more and more of the high-frequency information about the initial distribution specified by $p_0(\phi)$. This erasure might naïvely seem problematic since these flows appear to preclude the possibility of finding a non-trivial RG fixed point, importantly even in more general settings beyond ϕ^4 theory. However, the key is that to access a fixed point, we need to judiciously rescale the fields ϕ and couplings in a t -dependent manner, as discussed in the setting of a discrete-time flow in (4.23), (4.24). In other words, we need to appropriately ‘zoom in’ on the interesting fluctuations in our distribution, or else they will be lost to us.

More explicitly, suppose we have a distribution $p_{t, \{\lambda_i(t)\}}(\phi)$ where $\{\lambda_i(t)\}$ are the couplings in the action after an amount of RG flow t . To be maximally explicit, let us reprise our notation from Section 4.5 with some slight modifications. Suppose $p_{t, \{\lambda_i(t)\}}(\phi)$ is well approximated by

$$p_{t, \{\lambda_i(t)\}}(\phi) \approx \frac{1}{Z_t} \exp\left(-\left(\frac{L}{N}\right)^d \sum_i \lambda_i(t) M_{i,N}[\phi]\right) \quad (5.44)$$

where as before the $M_{i,N}[\phi]$'s given by

$$M_{i,N}[\phi] = \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \prod_{k=0}^{k_{\max,i}} (\Delta^k \phi(\mathbf{n}))^{q_{i,k}} \quad (5.45)$$

such that each $q_{i,k} \in \mathbb{Z}_{\geq 0}$. It is prudent for us to define

$$\delta_i := 2 \sum_{k=0}^{k_{\max,i}} k q_{i,k} \quad (5.46)$$

which counts the powers of $\epsilon = L/N$ appearing in $M_{i,N}[\phi]$ due to the Laplacians. Moreover, suppose that t corresponds to an effective cutoff scale Λ_t , which depends on the RG scheme employed. Then at some time t_* , the couplings $\{\lambda_i(t_*)\}$ define an (approximate) RG fixed point if for $t' \geq 0$ there exists a monotonic function $b_{t'}$ with $b_0 = 1$ such that

$$p_{t_*+t', \{\lambda_i(t_*+t')\}}(\phi) \prod_{\mathbf{n} \in \mathbb{Z}_N^d} d\phi(\mathbf{n}) \approx p_{t_*, \{(\frac{\Lambda_{t_*+t'}}{\Lambda_{t_*}})^{d+\delta_i} \lambda_i(t_*)\}}(b_{t'}\phi) \prod_{\mathbf{n} \in \mathbb{Z}_N^d} b_{t'} d\phi(\mathbf{n}). \quad (5.47)$$

An equivalent formulation is that there is a scaling function c_t so that

$$p_{t, \{(\frac{\Lambda_t}{\Lambda_0})^{d+\delta_i} \lambda_i(t)\}}(c_t\phi) \prod_{\mathbf{n} \in \mathbb{Z}_N^d} c_t d\phi(\mathbf{n}) \longrightarrow p_{\{\hat{\lambda}_i\}}^{\text{fixed pt}}(\phi) \prod_{\mathbf{n} \in \mathbb{Z}_N^d} d\phi(\mathbf{n}) \quad (5.48)$$

as t becomes large (but not so large that $\Lambda_t/\Lambda_0 \approx 0$), where $\{\hat{\lambda}_i\}$ is a fixed set of couplings.

We can easily accommodate for such a c_t in the flow equations by modifying the seed action $\widehat{S}_t(\tilde{\phi})$ as

$$\widehat{S}_t^{\text{new}}(\tilde{\phi}) := \widehat{S}_t(\tilde{\phi}) - \partial_t \log(c_t) \sum_{\mathbf{p} \in \mathbb{Z}_N^d} \frac{1}{\widetilde{B}_t(|\mathbf{p}|)} \tilde{\phi}(\mathbf{p}) \tilde{\phi}(-\mathbf{p}). \quad (5.49)$$

For instance, the SDE formulation of the Carosso scheme becomes

$$\partial_t \phi_t(\mathbf{n}) = (\Delta + \partial_t \log(c_t)) \phi(\mathbf{n}) + \eta_t(\mathbf{n}), \quad \phi_0(\mathbf{n}) = \phi(\mathbf{n}), \quad (5.50)$$

where the distribution of η_t is unchanged. To additionally implement the $\frac{\Lambda_t}{\Lambda_0}$ rescalings of the couplings $\lambda_i(t)$ (which is equivalent to a rescaling of the position or momenta), we can augment (5.50) by (see e.g. [14])

$$\partial_t \phi_t(\mathbf{n}) = (\Delta + \partial_t \log(c_t)) \phi(\mathbf{n}) - \partial_t \log(\Lambda_t) \mathbf{n} \cdot \mathbf{D}\phi(\mathbf{n}) + \eta_t(\mathbf{n}), \quad \phi_0(\mathbf{n}) = \phi(\mathbf{n}), \quad (5.51)$$

where

$$\mathbf{D}\phi(\mathbf{n}) := \left(\frac{1}{(\frac{L}{N})} (\phi(\mathbf{n} + \mathbf{e}_1) - \phi(\mathbf{n})), \frac{1}{(\frac{L}{N})} (\phi(\mathbf{n} + \mathbf{e}_2) - \phi(\mathbf{n})), \dots, \frac{1}{(\frac{L}{N})} (\phi(\mathbf{n} + \mathbf{e}_d) - \phi(\mathbf{n})) \right), \quad (5.52)$$

and again the distribution of η_t is unchanged.

We emphasize that the b_t or c_t required to access a fixed point is scheme-dependent; they will be different for e.g. the Carosso scheme versus the Polchinski scheme. The simplest way to see this is that the Carosso scheme and Polchinski scheme can suppress high-momentum modes at different rates.

A key problem is that, unless we have detailed analytic control over some combination of our RG scheme and our RG fixed point of interest, it is difficult to write down a suitable b_t or c_t a priori. This being said, if we suspect that a $p_{t, \{\lambda_i(t)\}}(\phi)$ is nearby a fixed point, we can use heuristics to search for an appropriate scaling function b_t , which if found would corroborate the presence of an RG fixed point. For the some, assume (somewhat unrealistically) that we have access to the explicit form of $p_{t, \{\lambda_i(t)\}}$. Then a straightforward heuristic is to fix a small time $\delta t' > 0$ and find a constant b such that

$$\mathbb{E}_{\phi \sim p_{t_* + \delta t', \{\lambda_i(t_* + \delta t')\}}(\phi)}[\phi(\mathbf{n})\phi(\mathbf{m})] \approx \frac{1}{b^2} \mathbb{E}_{\phi \sim p_{t, \{(\frac{\Lambda_{t_* + \delta t'}}{\Lambda_{t_*}})^{d + \delta_i} \lambda_i(t)\}}(\phi)}[\phi(\mathbf{n})\phi(\mathbf{m})], \quad (5.53)$$

where both sides are functions of $|\mathbf{n} - \mathbf{m}|_{\mathbb{Z}_N^d}$ due to translation-invariance. If such a b exists, then it is reasonable to guess that $b = b_{\delta t'}$. One can also look at higher-point analogs of (5.53).

In more realistic settings, we do not have direct access to the explicit form of $p_{t, \{\lambda_i(t)\}}$. Then it is standard to rely on *order parameters*, which are certain combinations of correlation functions of $p_{t, \{\lambda_i(t)\}}(\phi)$ (with no rescalings of the couplings $\{\lambda_i(t)\}$ or field ϕ) which diagnose the presence of a critical point. One drawback is that order parameters must be tailored to the critical point of interest; a priori, if we do not know anything about the critical point that an RG flow may land on, then it is unclear how to find an order parameter that detects it. For a review of order parameters in standard statistical field theories, see [72].

Now suppose we have two different Wegner-Morris RG flow schemes for the same initial $p_0(\phi)$; we denote their respective flows by $p_{t, \{\hat{\lambda}_i^{(1)}\}}^{(1)}(\phi)$ and $p_{t, \{\hat{\lambda}_i^{(2)}\}}^{(2)}(\phi)$, where $\{\hat{\lambda}_i^{(1)}\}$ denotes the fixed-point couplings for the first flow, and $\{\hat{\lambda}_i^{(2)}\}$ denotes the fixed-point couplings for the second flow. We do not in general expect $\{\hat{\lambda}_i^{(1)}\} \approx \{\hat{\lambda}_i^{(2)}\}$, nor do we expect $p_{t, \{\hat{\lambda}_i^{(1)}\}}^{(1)}(\phi) \approx p_{t, \{\hat{\lambda}_i^{(2)}\}}^{(2)}(\phi)$. However, universality of RG fixed point suggests that the two fixed points reached by two distinct RG flow schemes on the same initial $p_0(\phi)$ are related by a rescaling, namely that there is a constant c such that

$$p_{t, \{\hat{\lambda}_i^{(1)}\}}^{(1)}(\phi) \prod_{\mathbf{n} \in \mathbb{Z}_N^d} d\phi(\mathbf{n}) \approx p_{t, \{\hat{\lambda}_i^{(2)}\}}^{(2)}(c\phi) \prod_{\mathbf{n} \in \mathbb{Z}_N^d} c d\phi(\mathbf{n}). \quad (5.54)$$

The constant c can be ascertained by comparing the second moments of each distribution akin to our heuristic algorithm in (5.53); in particular:

$$\mathbb{E}_{\phi \sim p_{t, \{\hat{\lambda}_i^{(1)}\}}^{(1)}(\phi)}[\phi(\mathbf{n})\phi(\mathbf{m})] \approx \frac{1}{c^2} \mathbb{E}_{\phi \sim p_{t, \{\hat{\lambda}_i^{(2)}\}}^{(2)}(\phi)}[\phi(\mathbf{n})\phi(\mathbf{m})]. \quad (5.55)$$

Having discovered this c , which amounts to a choice of normalization of the ϕ field, we can now

explore the same aspects of the RG fixed point with each scheme in a manner such that the results will (approximately) agree.

In other applications of RG on the lattice, one considers theories which are well-approximated by a finite number of terms in the action at long distances (i.e. a *renormalizable* EFT). Let the couplings associated to those terms be called $\{\lambda_i(t)\}$. One desires to pick initial values of those couplings $\{\lambda_i(0)\}$ for the theory at short distances (in the UV) so that they flow to a desired set of couplings $\{\lambda_i(T)\}$ at the long distance scale corresponding to T (in the IR); all other possible couplings not in the set can be initially set to zero since they will be suppressed at long distances anyway. More specifically, suppose that for a particular renormalization group scheme, all momentum modes with $|\widehat{\mathbf{p}}_i| \geq \Lambda_0 f(T)$ are suppressed, and all momentum modes with $|\widehat{\mathbf{p}}_i| \leq \Lambda_0 f(T)$ are unsuppressed. Let us say that we have a physical system for which we have measured (some of) the effective couplings at the momentum scale corresponding to $\Lambda_0 f(T)$, for a fixed T . Then, to match this to lattice data, we would like to tune the $\{\lambda_i(0)\}$ so that they flow to the right couplings at $t = T$. Having identified the appropriate $\{\lambda_i(0)\}$, our lattice model can now make predictions about the values of the couplings at $t \neq T$, for instance correlation functions at momentum scales other than $\Lambda_0 f(T)$.

Some of the ideas and techniques for RG in this section also apply to models outside the purview of Effective Field Theory, such as image models. Perhaps a useful insight in this more general context is as follows. We know that RG flows are sensitive to rescaling ϕ along the flow, especially if want to access interesting correlations. As such, it could be prudent to consider the pushforward of $p_t(\phi)$ under rescalings $\phi \rightarrow b\phi$ by a judicious constant b which may allow one to ‘zoom in’ on the interesting correlations. Such a b could be identified in the following way. Let $F_b(\phi) = b\phi$, and define $\bar{\phi}(\mathbf{n}) := \mathbb{E}_{\phi \sim F_b * p_t(\phi)}$. Then we may desire to pick a b such that e.g.

$$\Phi(\mathbf{n}) := \sum_{\mathbf{m}} \mathbb{E}_{\phi \sim F_b * p_t(\phi)} [(\phi(\mathbf{n}) - \bar{\phi}(\mathbf{n}))(\phi(\mathbf{m}) - \bar{\phi}(\mathbf{m}))] \quad (5.56)$$

is ≈ 1 for \mathbf{n} near the center of the lattice (if it does not have periodic boundary conditions, as in the case of an image). In other words, we let b set the scale of fluctuations in the model.

Another general feature which may be interesting in general models is to choose a $\widehat{S}_t(\phi)$ which contains terms non-quadratic in ϕ , for instance quartic. This was discussed in [124] (see also comments in [30]), and may ameliorate the ambiguity of the overall scale of ϕ near fixed points; that is, the nonlinearity of the RG flow in some instances may pick out a particular scaling. At the present time, this avenue of non-quadratic $\widehat{S}_t(\phi)$ ’s has been less explored.

In sum, we emphasize that the Wegner-Morris-type RG schemes offer broad flexibility for RG flows with desirable properties in momentum space. Their usage, in physical applications or more generally, needs to be augmented by novel physical intuitions coming from RG theory in order to robustly access desired correlations in the ensuing $p_t(\phi)$ ’s.

6 Finding ground states of quantum field theories

6.1 Review of difficulties with variational methods in quantum field theory

A fundamental problem in the study of quantum mechanics is finding the ground state of a system. At a formal level, quantum systems are in part described by a Hermitian operator H called the *Hamiltonian*, which has bounded spectrum from below. Suppose, for simplicity, that there is a unique ground state of H , i.e. the eigenvector $|\Psi_{\text{gs}}\rangle$ of H with the smallest eigenvalue E_{gs} . The eigenvalue E_{gs} is called the *ground state energy*. As such, there is a variational principle that recovers the ground state, namely

$$|\Psi_{\text{gs}}\rangle = \arg \min_{|\Psi\rangle} \frac{\langle \Psi | H | \Psi \rangle}{\langle \Psi | \Psi \rangle}. \quad (6.1)$$

This is called the *Rayleigh-Ritz* variational principle. (Above, $\langle \Psi |$ is the Hermitian conjugate of the vector $|\Psi\rangle$.)

The Rayleigh-Ritz variational principle is often used in the following manner. Suppose we parameterize some submanifold of the space of $|\Psi\rangle$'s with some parameters θ . The $|\Psi\rangle$ corresponding to a fixed θ is denoted by $|\Psi^\theta\rangle$. Then a proxy for the minimization in (6.1) is

$$\tilde{\theta} := \arg \min_{\theta} \frac{\langle \Psi^\theta | H | \Psi^\theta \rangle}{\langle \Psi^\theta | \Psi^\theta \rangle}, \quad (6.2)$$

where $|\Psi^{\tilde{\theta}}\rangle$ is an approximation to the true ground state $|\Psi_{\text{gs}}\rangle$, and $\frac{\langle \Psi^{\tilde{\theta}} | H | \Psi^{\tilde{\theta}} \rangle}{\langle \Psi^{\tilde{\theta}} | \Psi^{\tilde{\theta}} \rangle}$ is an approximation to (and in fact an upper bound for) the true ground state energy E_{gs} .

In practical applications, one attempts to formulate a judicious parameterization $|\Psi^\theta\rangle$ so that, for a given class of Hamiltonians H , the optimal $|\Psi^{\tilde{\theta}}\rangle$ should be close to the true ground state $|\Psi_{\text{gs}}\rangle$. However, it is most often intractable to prove that a particular parameterization of convenience can be optimized so that the resulting state is in close proximity to a desired ground state. As such, variational methods are often used as heuristics, and compared with other methods and forms of data including experiments coming from natural systems. A highly-successful specialization of these methodologies is Density Functional Theory (DFT), which is a workhorse of calculations in quantum chemistry [125–127].

Here we will focus our attention on the setting of quantum field theory, where the application of variational methods is both highly desirable and difficult. To understand the source of the essential difficulty, we recall some basic facts about quantum field theory via a standard example. Earlier in this paper, we considered scalar ϕ^4 theory as a *Euclidean field theory*, e.g. as a *statistical field theory*. Now we consider its quantum mechanical counterpart, and focus on the lattice setting in particular.

As usual, we consider a d -dimensional lattice in the hypercube with side length L , and lattice sites $\frac{L}{N} \mathbf{n}$ for $\mathbf{n} \in \{0, 1, \dots, N\}^d$. We periodically identify opposing sides of the hypercube so that it becomes a torus, and thus $\mathbf{n} \in \mathbb{Z}_N^d$. To each lattice site we associate the Hilbert space $L^2(\mathbb{R})$, so

that the total Hilbert space is $\mathcal{H} \simeq \bigotimes_{\mathbf{n} \in \mathbb{Z}_N^d} \mathcal{H}_{\mathbf{n}} \simeq (L^2(\mathbb{R}))^{\otimes (\frac{L}{N})^d}$. Here each $\mathcal{H}_{\mathbf{n}} \simeq L^2(\mathbb{R})$ is a space of L^2 functions, which we associate with the variable $x_{\mathbf{n}}$. So, for instance, a function on $\mathcal{H}_{\mathbf{n}}$ will be denoted by $f(x_{\mathbf{n}})$, a derivative of such a function is denoted with a $\frac{\partial}{\partial x_{\mathbf{n}}}$, and so on. Now the Hamiltonian H operator acts on \mathcal{H} , where the Hamiltonian is given by

$$H := \left(\frac{L}{N}\right)^d \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \left(-\frac{1}{2} \frac{\partial^2}{\partial x_{\mathbf{n}}^2} + \frac{1}{2 \left(\frac{L}{N}\right)^2} \sum_{i=1}^d (x_{\mathbf{n}+\mathbf{e}_i} - x_{\mathbf{n}})^2 + \frac{1}{2} m^2 x_{\mathbf{n}}^2 + \lambda x_{\mathbf{n}}^4 \right). \quad (6.3)$$

We can view a state $\Psi(\{x_{\mathbf{n}}\})$ in the Hilbert space \mathcal{H} as living on a d -dimensional spatial lattice, and the Schrödinger equation tells us that it evolves in time t as $e^{-iHt}\Psi(\{x_{\mathbf{n}}\})$. For physically-relevant systems, we essentially always have $d = 1, 2, 3$.

With our Hilbert space and Hamiltonian at hand, it is useful to understand the dimensionality of said Hilbert space. Each $L^2(\mathbb{R})$ tensor factor is infinite-dimensional, but in practice we can truncate each into a k -dimensional vector space for some suitable k . With this truncation, the total Hilbert space dimension is $k^{\left(\frac{L}{N}\right)^d}$, which is enormous. This has grave implications for applying (6.1) to our Hamiltonian H : if we parameterize an arbitrary state in the Hilbert space, it requires $k^{\left(\frac{L}{N}\right)^d}$ numbers to specify, and so the minimization would be over $k^{\left(\frac{L}{N}\right)^d}$ parameters. This is clearly intractable even for modest system sizes. As such, to apply the Rayleigh-Ritz variational principle to the setting of quantum field theory, it is imperative to find an efficient parameterization $\Psi^{\theta}(\{x_{\mathbf{n}}\})$ involving a *sub-exponential* number of parameters θ , such that a minimization like (6.2) over θ will lead to a good approximation for the ground state of the quantum field theory. However, there are interesting technical obstructions to finding such an efficient parameterization.

As pointed out by Feynman in [35], a central difficulty in the setting of quantum field theory is that the minimization in the Rayleigh-Ritz variational principle is overly sensitive to changes in the wavefunction $\Psi^{\theta}(\{x_{\mathbf{n}}\})$ at the lattice scale. For example, this sensitivity can be gleaned by examination of the $\left(\frac{L}{N}\right)^d \frac{1}{2 \left(\frac{L}{N}\right)^2} \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \sum_{i=1}^d (x_{\mathbf{n}+\mathbf{e}_i} - x_{\mathbf{n}})^2$ part of the Hamiltonian (6.3). In physical terms, changes of the parameters θ that affect $\Psi^{\theta}(\{x_{\mathbf{n}}\})$ at the lattice scale have an outsized energetic cost, whereas changes of θ that affect longer-range correlations have comparatively smaller energetic cost. This presents a problem because, in practice, the medium-range and longer-range correlations are essential for comparing to physical experiments, whereas the short-range correlations tied to the lattice scale can be viewed as non-universal artifacts. For instance, if our lattice theory is viewed as an approximation to a continuum theory at long distances, then we desire that the lattice nature of our theory (often implemented for numerical convenience) does not hold hostage the accuracy of long-range correlations when a variational optimization is performed.

One strategy to ameliorate Feynman's roadblock is to choose hierarchical ansätze for the θ -dependence of $\Psi^{\theta}(\{x_{\mathbf{n}}\})$ so that θ more equitably controls correlations in the wavefunction across all distance scales. At a heuristic level, we can imagine performing RG flow in reverse: we start with a fiducial state which parameterizes long-distance correlations and we allot some part of the θ -parameters to describe this state; then we perform some kind of reverse-RG flow to build

up shorter-distance correlations on top of the longer-distances ones, wherein the new shorter-distance correlations are described by another part of the θ parameters. This is then repeated over multiple rounds until we have a state $\Psi^\theta(\{x_{\mathbf{n}}\})$ for which the parameters θ equitably parameterize correlations across all scales. This strategy has been concretely implemented in the context of *tensor networks* [128], which have been remarkably successful for providing tractable and viable ansätze for the Rayleigh-Ritz variational principle in the context of field theories in one spatial dimension [36–38] (i.e. $d = 1$). However, the setting of two and three spatial dimensions (i.e. $d = 2, 3$) remains mostly out of reach with current tensor network methods (see e.g. [129] for a discussion of difficulties in $d = 2$), and so new ansätze are required. In particular, existing tensor network ansätze in higher spatial dimensions are either not computationally practical, or if they are then they are not capacious enough to well-approximate the true desired ground state wavefunction. In the Subsections below we describe a class of potentially useful ansätze inspired by the reverse-RG logic explained here, and anticipate that this class may have great utility in describing ground states of field theories for all of $d = 1, 2, 3$.

6.2 Real-valuedness of the ground state wavefunction

Here we explain a useful fact about the ground state wavefunction $\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})$ of a Hamiltonian H with a unique ground state. The fact will be useful in our variational algorithm described in the next Subsection.

Suppose that $\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})$ is the ground state wavefunction of H , and that it is complex-valued. If the energy of the ground state is E_{gs} , then we have

$$(H - E_{\text{gs}}) \text{Re}\{\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})\} + i(H - E_{\text{gs}}) \text{Im}\{\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})\} = 0, \quad (6.4)$$

and due to the Hermiticity of H we find that $(H - E_{\text{gs}}) \text{Re}\{\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})\} = 0$ and $(H - E_{\text{gs}}) \text{Im}\{\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})\} = 0$ individually. But if H has a unique ground state, then $\text{Re}\{\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})\}$ must be proportional to $\text{Im}\{\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})\}$. But this means that, if $\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})$ is L^2 -normalized, we must have

$$\Psi_{\text{gs}}(\{x_{\mathbf{n}}\}) = e^{i\varphi} \text{Re}\{\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})\} \quad (6.5)$$

for some phase φ . But since wavefunctions are only defined up to a global phase (since a global phase $e^{i\varphi}$ does not affect any measurable quantity in a quantum-mechanical theory), it follows that we can take $\varphi = 0$ so that $\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})$ is purely real-valued.

In summary, if we have a quantum field theory with a unique ground state wavefunction, then it can be taken to be real. This holds for the example of (6.3) above, and for many other examples.

6.3 Learning ground states of QFTs with diffusion models

In this Subsection we develop a variational algorithm for learning ground states of quantum lattice systems such as (6.3). The basic idea is to leverage the stochastic formalism of the Exact Renormalization Group from Section 5, and in particular to variationally perform ERG in reverse on a

fiducial state to build up correlations at progressively smaller distance scales in such a way that we ultimately arrive at a good approximation to our desired ground state wavefunction. Before describing our algorithm, let us develop some further notation.

Let us reprise the Hamiltonian in (6.3), writing it in a slightly more compact form as

$$H = \left(\frac{L}{N}\right)^d \left(-\frac{1}{2} \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \frac{\partial^2}{\partial x_{\mathbf{n}}^2} + \mathcal{F}(\{x_{\mathbf{n}}\}) \right). \quad (6.6)$$

We previously denoted the ground state wavefunction by the notation $\Psi_{\text{gs}}(\{x_{\mathbf{n}}\})$, which is a function of N^d variables. But now let us define a function

$$\phi : \mathbb{Z}_N^d \longrightarrow \mathbb{R}, \quad \phi(\mathbf{n}) = x_{\mathbf{n}}, \quad (6.7)$$

so that we have

$$\frac{\partial}{\partial \phi(\mathbf{n})} \longleftrightarrow \frac{\partial}{\partial x_{\mathbf{n}}}, \quad d\phi \longleftrightarrow \prod_{\mathbf{n} \in \mathbb{Z}_N^d} dx_{\mathbf{n}}. \quad (6.8)$$

With these scalar field notations, we can rewrite the Hamiltonian (6.6) as

$$H = \left(\frac{L}{N}\right)^d \left(-\frac{1}{2} \Delta_{\phi} + \mathcal{F}(\phi) \right), \quad (6.9)$$

where $\Delta_{\phi} := \nabla_{\phi} \cdot \nabla_{\phi} = \sum_{\mathbf{n} \in \mathbb{Z}_N^d} \frac{\partial}{\partial \phi(\mathbf{n})}$, and we write the ground state wavefunction as $\Psi_{\text{gs}}(\phi)$. We will similarly write our variational wavefunction as $\Psi^{\theta}(\phi)$.

For our purposes, let us write the variational wavefunction as

$$\Psi_t^{\theta}(\phi) = \frac{1}{\sqrt{Z_t^{\theta}}} e^{-S_t^{\theta}(\phi)/2} \quad (6.10)$$

where $\Psi_{t=0}^{\theta}(\phi)$ will ultimately correspond to our best variational approximation to the ground state of H . Above, $S_t^{\theta}(\phi)$ is not an action as in the statistical field theory context, but is simply a convenient family of functions. Note, however, that $S_t^{\theta}(\phi)$ can be chosen to be real on account of our discussion in Section 6.2. The factor Z_t^{θ} is a normalizing constant that we will not need explicitly, but such that the wavefunction is normalized as

$$\int d\phi |\Psi_t^{\theta}(\phi)|^2 = \frac{1}{Z_t^{\theta}} \int d\phi e^{-S_t^{\theta}(\phi)} = 1. \quad (6.11)$$

It will often be convenient to write

$$p_t^{\theta}(\phi) := |\Psi_t^{\theta}(\phi)|^2. \quad (6.12)$$

The energy of the state $\Psi_t^{\theta}(\phi)$ can be written as

$$\frac{1}{Z_t^{\theta}} \int d\phi e^{-S_t^{\theta}(\phi)/2} H e^{-S_t^{\theta}(\phi)/2} = \mathbb{E}_{\phi \sim p_t^{\theta}(\phi)} \left[-\frac{1}{2} \Delta_{\phi} S_t^{\theta} + \frac{1}{4} (\nabla_{\phi} S_t^{\theta})^2 + \mathcal{F}(\phi) \right]. \quad (6.13)$$

One can compute the θ -gradient of this function using the formula $\nabla_\theta e^{-S_t^\theta/2} = e^{-S_t^\theta/2}(-\nabla_\theta S_t^\theta/2)$. Unfortunately, this means that we need estimates for S_t^θ to compute the gradients of the function; fortunately, such estimates are provided by using normalizing flows and Neural ODE methods. We present one possible approach below. Recall that the probability flow ODE associated to (5.25) is

$$\frac{d\phi}{dt} = \Delta\phi - \frac{1}{2}\sigma^2 s_t^\theta(\phi) =: b_\theta(\phi, t), \quad (6.14)$$

with $\epsilon = L/N$, and $s_t^\theta := \nabla_\phi \log p_t^\theta$. Supposing that $p_\infty(\phi)$ is some distribution of our choice which is easy to sample from, we have the following variational algorithm for determining $\Psi_t^\theta(\phi)$:

initialize $\theta = (\theta^1, \dots, \theta^M)$, $i = 0$

if $i < i_{\max}$ **do**

sample $\phi'_T \leftarrow p_T^\theta = p_\infty$

compute ϕ'_0 as well as $\delta\theta(0)$ via an ODE solver, namely, by solving (6.14) as well as the system

from $t = T$ to $t = 0$ with initial condition $x = \phi'_T$, $\delta_{\theta^k} x(T) = 0$, $\delta\theta^k(T) = 0$ for all k :

$$\frac{dx}{dt} = b_\theta(x, t), \quad \frac{d(\delta_{\theta^k} x)}{dt} = \frac{\partial}{\partial\theta^k} b_\theta(x, t), \quad \frac{d(\delta\theta^k)}{dt} = \nabla_x \cdot \left(\frac{\partial}{\partial\theta^k} b_\theta(x, t) \right) + \sum_\ell \sum_m b^\ell(x, t) \delta_{\theta^m} x,$$

set $\bar{\theta}_0 \leftarrow (\delta\theta^1(0), \dots, \delta\theta^M(0))$

set $\theta_0 \leftarrow \bar{\theta}_0 (-\nabla_\phi \cdot s_0^\theta(\phi)/2 + |s_0^\theta(\phi)|^2 + \mathcal{F}(\phi))|_{\phi=\phi'_0} + (-\nabla_\theta \nabla_\phi \cdot s_0^\theta(\phi) - \nabla_\theta |s_0^\theta(\phi)|^2)|_{\phi=\phi'_0}$

for $j = 1, \dots, \tau$

sample $t_j \sim [0, T]$

sample $\phi_j \leftarrow p_{t_j}^\theta$ using (5.19)

compute $\nabla_\theta \log p_{t_j}^\theta(\phi_j)$ by calling an ODE solver, namely:

(i) solve $dx/dt = b_\theta(x, t)$ from $t = t_j$ to $t = T$ with initial condition $x_t = \phi_t$

(ii) set $x_T \leftarrow x(T)$

solve the following system from $t = T$ to $t = t_j$ for $k = 1, \dots, M$:

$$\frac{dx}{dt} = b_\theta(x, t), \quad \frac{d(\delta_{\theta^k} x)}{dt} = \frac{\partial}{\partial\theta^k} b_\theta(x, t), \quad \frac{d(\delta\theta^k)}{dt} = \nabla_x \cdot \left(\frac{\partial}{\partial\theta^k} b_\theta(x, t) \right) + \sum_\ell \sum_m b^\ell(x, t) \delta_{\theta^m} x,$$

with initial conditions $\delta_{\theta^k} x(T) = 0$, $\delta\theta^k(T) = 0$ for all k where $\delta\theta^k(t) := \frac{\partial}{\partial\theta^k} \log p_t^\theta(x(t))$

set $\theta_j \leftarrow (\delta\theta^1(t_j), \dots, \delta\theta^M(t_j))$

end for

set $\theta \leftarrow \theta - \sum_{j=0}^{\tau} \theta_j$

end if

return θ

Note that while the algorithm provides us with $\Psi_{t=0}^\theta(\phi)$ as an approximation to the ground state $\Psi_{\text{gs}}(\phi)$, we also end up with $\Psi_t^\theta(\phi)$ for various values of t , which can be thought of as (an approximation to) the ground state RG-flowed by different amounts.

7 Numerically learning RG flows

7.1 Overview

We tested the methods of Section 5 on the simple examples of 2D lattice scalar field theories. Specifically, we tested that the normalizing-flow based method based on optimizing (5.22) can learn the RG flow of such theories under the Carosso scheme (5.9). We parameterized the normalizing flow b_t^θ of (5.8) using the neural network architecture proposed by [4]. We used a batch size of 64, namely we computed gradient updates for 64 independent samples treated as in the algorithm described below (5.22), and fed the averaged gradients into the Adam optimizer [130]. For all experiments, we used hyperparameters $b_1 = 0.8$, $b_2 = 0.9$ for the optimizer provided by the Optax package [131], and used an exponentially decaying learning rate with initial value 0.005 followed by 8000 transition steps and decay rate 0.1.

We show that this method learns the flows of basic physical quantities like the renormalized mass. We present our lattice field theory conventions, as well as the estimators we use for the relevant physical quantities, in Section 7.2. We plot the flows of these quantities using the Carosso and Polchinski SDEs in Section 7.3 to get a sense of the kinds of differences that are possible when using these two radically different schemes.

We found that it is difficult to learn the RG flows of field theories when the RG flow is defined using the Polchinski scheme. This is likely because the flow of the Polchinski scheme rapidly compresses the support of the distribution p_t to a very low-dimensional manifold, and the inverse flow has trouble spreading out this low-dimensional prior to the UV distribution over lattice fields, which is supported everywhere. These issues with the Polchinski scheme might be ameliorated by directly incorporating field rescalings and spatial/momentum rescalings into the flow equations, as discussed around (5.51).

With the Carosso scheme, we found that learning the RG flow was not possible unless a small “mass” term was added to modify the RG SDE to

$$\partial_t \phi_t(\mathbf{n}) = (\Delta - M^2)\phi(\mathbf{n}) + \eta_t(\mathbf{n}), \quad \phi_0(\mathbf{n}) = \phi(\mathbf{n}), \quad (7.1)$$

which the noise distribution unchanged. We previously discussed this modification around (5.20). Note that the unmodified Carosso scheme does not change the *mean value of the field*, i.e. the distribution over paths in field-space ϕ_t induced by (5.9) is invariant under $\phi_t \mapsto \phi_t + C$ for any constant C . As such, there is a 1-dimensional line of limiting distributions p_∞ ; in contrast, with the above modification in (7.1), there is only a single limiting distribution. We use $M = 1$ for the experiments below; we find that the qualitative features of all plots are not sensitive to the value of M so long as it is sufficiently small. We also make a similar M -modification to the Polchinski scheme. The corresponding transition kernels of the M -modified Carosso and Polchinski schemes are (5.20) and (5.41) with the $K_t(|\hat{\mathbf{p}}|^2)$ in (5.42).

7.2 Conventions and estimators

Let us recapitulate our conventions for scalar ϕ^4 theory in two dimensions. We have the probability distribution

$$p_0(\phi) = \frac{1}{Z} \exp \left(-L^2 \sum_{\mathbf{p} \in \mathbb{Z}_N^2} \frac{1}{2} \tilde{\phi}(\mathbf{p})(|\hat{\mathbf{p}}|^2 + m^2) \tilde{\phi}(-\mathbf{p}) - L^2 \sum_{\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3, \mathbf{p}_4 \in \mathbb{Z}_N^2} \frac{\lambda}{4!} \tilde{\phi}(\mathbf{p}_1) \tilde{\phi}(\mathbf{p}_2) \tilde{\phi}(\mathbf{p}_3) \tilde{\phi}(\mathbf{p}_4) \delta_{\mathbf{p}_1 + \mathbf{p}_2 + \mathbf{p}_3 + \mathbf{p}_4, \mathbf{0}} \right), \quad (7.2)$$

where m and λ are the initial bare mass and bare coupling, respectively. In our numerics below, we present statistical estimates for four quantities along the RG flow of a scalar ϕ^4 theory: (i) the renormalized mass m_R , (ii) the wave function renormalization Z , (iii) the renormalized coupling λ_R , and (iv) a diagnostic of field excursions $\max_{\mathbf{n}} |\phi(\mathbf{n})|$.

We will define the renormalized couplings m_R and λ_R in the standard way, which we review here. In the M -modified Carosso scheme (and the M -modified Polchinski scheme using the $K_t(|\hat{\mathbf{p}}|^2)$ in (5.42)), an effective cutoff scale corresponding to t is

$$\Lambda_t = \max \left\{ \min \left\{ \sqrt{\frac{1}{t} - M^2}, \Lambda_0 \right\}, 0 \right\}. \quad (7.3)$$

Defining $r_t := \frac{\Lambda_t}{\Lambda_0}$, we have

$$r_t = \frac{\Lambda_t}{\Lambda_0} := \max \left\{ \min \left\{ \sqrt{\frac{1}{\Lambda_0^2 t} - \frac{M^2}{\Lambda_0^2}}, 1 \right\}, 0 \right\}. \quad (7.4)$$

The Carosso scheme does not explicitly account for rescaling momenta by $\mathbf{p} \rightarrow r_t \mathbf{p}$ and fields by $\tilde{\phi} \rightarrow \tilde{\phi}/r_t^2$ (i.e. accounting for non-anomalous dimensions of $\tilde{\phi}$ in $d = 2$). However, these rescalings are necessary to provide the usual definitions of the renormalized couplings, so we need to include these rescalings by hand. We incorporate these rescalings below.

To estimate our renormalized couplings, we need to define a few correlation functions. The first is the momentum space 2-point function, namely

$$\tilde{G}_2^{\text{Carosso}}(\mathbf{p}_1, \mathbf{p}_2) := \langle \tilde{\phi}(\mathbf{p}_1) \tilde{\phi}(\mathbf{p}_2) \rangle. \quad (7.5)$$

The next one is the momentum space connected 4-point function, given by

$$\begin{aligned} \tilde{G}_{4, \text{conn}}^{\text{Carosso}}(\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3, \mathbf{p}_4) &:= \langle \tilde{\phi}(\mathbf{p}_1) \tilde{\phi}(\mathbf{p}_2) \tilde{\phi}(\mathbf{p}_3) \tilde{\phi}(\mathbf{p}_4) \rangle - \langle \tilde{\phi}(\mathbf{p}_1) \tilde{\phi}(\mathbf{p}_2) \rangle \langle \tilde{\phi}(\mathbf{p}_3) \tilde{\phi}(\mathbf{p}_4) \rangle \\ &\quad - \langle \tilde{\phi}(\mathbf{p}_1) \tilde{\phi}(\mathbf{p}_3) \rangle \langle \tilde{\phi}(\mathbf{p}_2) \tilde{\phi}(\mathbf{p}_4) \rangle - \langle \tilde{\phi}(\mathbf{p}_1) \tilde{\phi}(\mathbf{p}_4) \rangle \langle \tilde{\phi}(\mathbf{p}_2) \tilde{\phi}(\mathbf{p}_3) \rangle. \end{aligned} \quad (7.6)$$

For small $|\hat{\mathbf{p}}|$, the 2-point function goes as

$$\tilde{G}_2^{\text{Carosso}}(\mathbf{p}, -\mathbf{p}) \approx \frac{Z}{|\hat{\mathbf{p}}|^2 + r_t^2 m_R^2}. \quad (7.7)$$

Then we can formulate estimators for m_R and Z by

$$\frac{1}{m_R^2} = \xi^2 \approx \frac{r_t^2}{4L^2} \sum_{\mathbf{p} \in \{(1,0), (0,1), (-1,0), (0,-1)\}} \frac{1}{\frac{4N^2}{L^2} \sin^2(\pi/N)} \left(\frac{\tilde{G}_2^{\text{Carosso}}(\mathbf{0}, \mathbf{0})}{\tilde{G}_2^{\text{Carosso}}(\mathbf{p}, -\mathbf{p})} - 1 \right) \quad (7.8)$$

and

$$Z \approx r_t^2 m_R^2 \tilde{G}_2^{\text{Carosso}}(\mathbf{0}, \mathbf{0}) \quad (7.9)$$

where $\xi = 1/m_R$ is known as the correlation length. For the renormalized quartic coupling we use the standard estimator (see e.g. [132] or [133] and references therein)

$$\lambda_R \approx -r_t^2 (r_t m_R L)^2 \frac{\tilde{G}_{4, \text{conn}}^{\text{Carosso}}(\mathbf{0}, \mathbf{0}, \mathbf{0}, \mathbf{0})}{\tilde{G}_2^{\text{Carosso}}(\mathbf{0}, \mathbf{0})} \quad (7.10)$$

Finally, to estimate extreme field excursions, we will also estimate

$$\max_{\mathbf{n} \in \mathbb{Z}_N^2} |\phi(\mathbf{n})| \quad (7.11)$$

7.3 The Carosso and Polchinski flows

In Figure 6, we show plots of the flows of some of the above estimators using the Carosso and the Polchinski schemes (with the M -modifications discussed previously). We see that while UV samples look completely different when flowed using the Carosso scheme or the Polchinski scheme, the resulting flows of estimators of physical quantities are qualitatively similar. In particular, both the mass and the quartic couplings flow at similar rates; the quartic coupling (multiplied by r_t^2 , for ease of visualization) goes to zero slightly faster in the Polchinski scheme than in the Carosso scheme, but only by a factor of two or so. In the Carosso scheme, $r_t^2 \lambda_R$ ultimately goes to zero for larger values of t than we show in the plots.

Thus, by picking a convenient exact renormalization group scheme, we are still able to get physically meaningful results while dramatically improving the numerical stability of the training process for the sampler. This highlights the importance of tuning the RG scheme, just as the corresponding tuning of the noising process is necessary to ensure good performance of generative models of images.

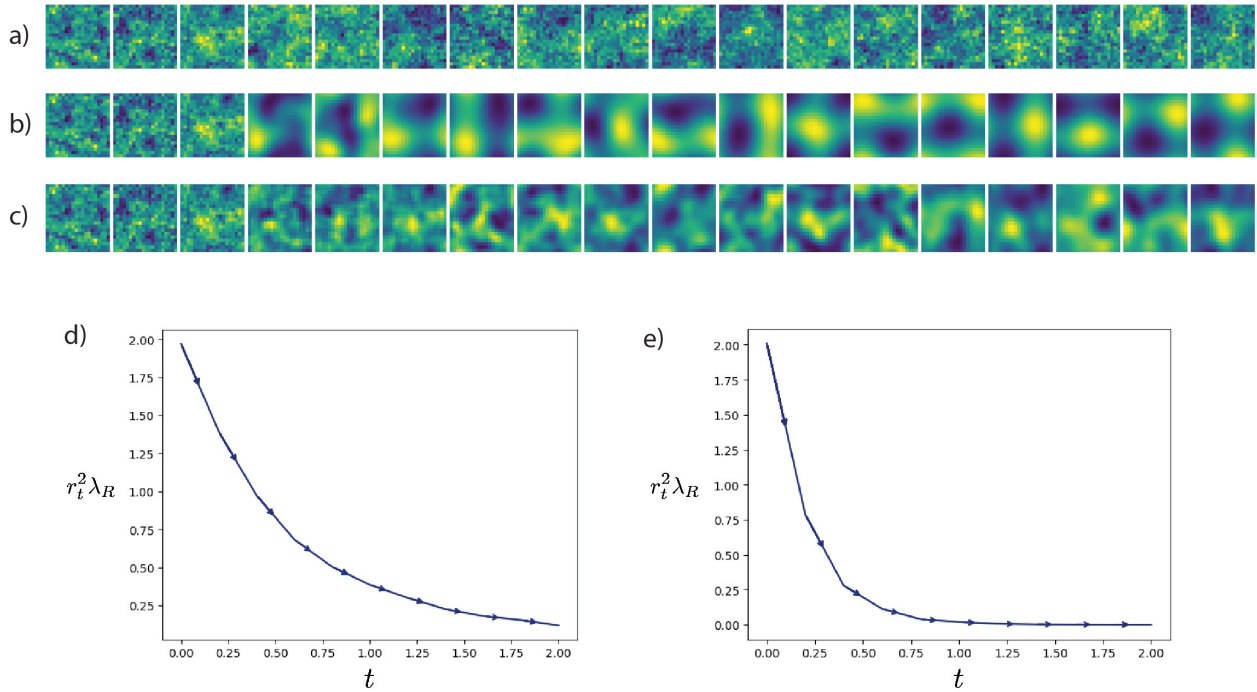


Figure 6: *ERG schemes can be qualitatively and numerically different but physically equivalent.* **Top:** Flows of samples from lattice ϕ^4 theory in two dimensions (with $L = 1$, $N = 20$, $m^2 = -2$, $\lambda = 0.0017$) via (a) the Carosso scheme and (b) the Polchinski scheme. (c) The third row slows down the Polchinski flow, corresponding to the first few frames of (b). Note how the Polchinski flow cuts off high frequencies extremely rapidly, making it difficult to train models that learn the Polchinski flow. **Bottom:** Plot of renormalized quartic coupling (rescaled by r_t^2 for clarity) under the Carosso flow (d) and the Polchinski flow (e), over the entire scale range of flows from the top of Figure in (a) and (b). Flows of couplings are at similar rates to within a small factor, despite the underlying samples being very different in nature, as expected from the ERG formalism. Note that $r_t^2 \lambda_R$ goes to zero for large t in each scheme, although in the Carosso scheme $t = 2$ is not long enough to see $r_t^2 \lambda_R$ reach near zero.

7.4 Numerics

We also compared the normalizing flows trained using the reverse KL objective, as in (5.5) (following [2]), with the flows trained using our objective (5.22). In Figures 7 and 8, we plot the estimated flows of the physical quantities described above, as well as the flows of the same quantities estimated from samples drawn directly from the learned normalizing flows. In all cases, we trained the reverse KL model and the model based on (5.22) for the same number of gradient steps (1000). To compute estimates of the flows of the physical quantities, we performed sampling from the UV distribution using the NUTS sampler [134], as implemented in the Blackjax package [135], with an

initial warm-up and tuning segment of 4000 time-steps using the “window adaptation” method (the standard adaptation method used in Stan [136]) implemented in [135]. We then evolved forwards in t with the Carosso scheme using the kernels derived in this paper.

We found that according to a variety of estimates of the flows of physically-important quantities, our objective caused the learned normalizing flow to more accurately reflect the physical behavior of the RG scheme, while the method based on optimizing the reverse KL divergence did not reliably lead to any agreement with the RG flow. However, the reverse KL divergence was often able to fit the relevant estimators for the UV distribution, i.e. at RG time $t = 0$. Moreover, we found that this difference in learned flows, which is evident from the flows of physical quantities, cannot be detected by observing sampled field configurations by eye. This may come as no surprise to those who study lattice field theory, but in the context of machine learning, visual comparisons often play a significant role in identifying a successful model. We found that estimators of physical quantities were very sensitive to the distribution we tried to learn, and that looking at the behavior of these estimators under RG flow allowed us to quickly identify bugs in our training code, since we knew what qualitative behavior to expect from the physics of the ϕ^4 model.

We also note that in all of the experiments shown, we used the limiting distribution for the Carosso scheme as the prior distribution for the sampler, even on the samplers trained using the reverse-KL objective. While this is necessary for the consistency of the forwards-KL samplers, it is not required for the reverse-KL samplers. In fact, it is more common [2, 4] to use a white noise distribution as a prior for a reverse-KL trained sampler. If we compared our flows (which have a Carosso-type prior) to flows generated by reverse-KL-trained samplers with a white noise prior, we would find that the reverse-KL-trained samplers would be significantly worse at learning the RG flow of the true field theory than the reverse-KL-trained flows illustrated in Figures 7 and 8; thus our method would appear to perform significantly better. Instead, we opted to make a more challenging comparison (with respect to the problem of learning the RG flow of the field theory) by comparing our method with reverse-KL samplers trained with a Carosso-type prior. The Carosso-type prior encourages the reverse-KL-trained normalizing flow to match the ‘true’ RG flows because we are effectively forcing the flows induced by the samplers to have the correct initial and final points; since the learned flows are continuous, they will at least qualitatively look like the ‘true’ RG flows. Nevertheless, even with our more careful comparison, we still show a significant improvement of our method over previous ones in the context of matching to RG flows.

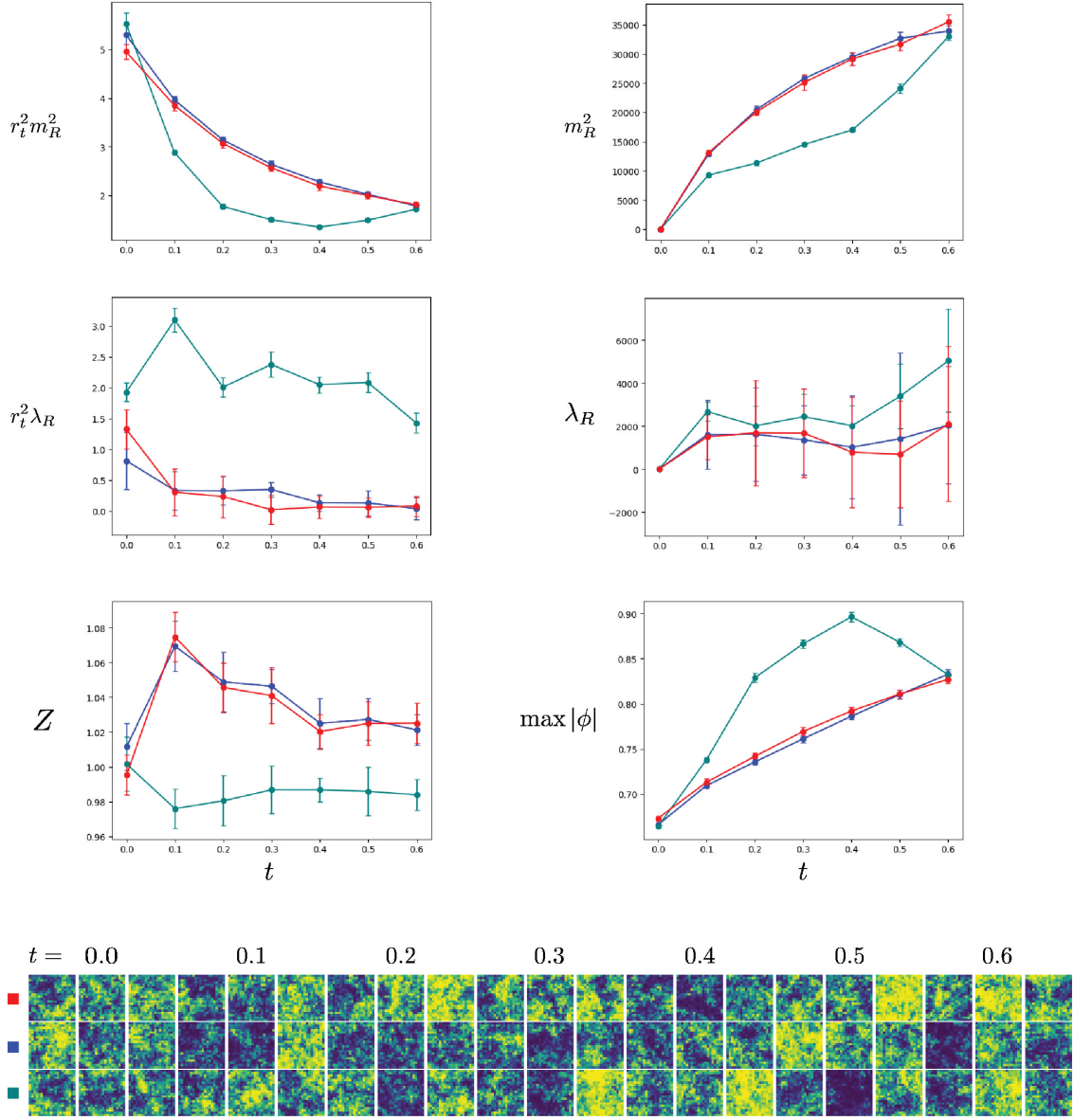


Figure 7: *Learning RG Flow of ϕ^4 theory with bare parameters $L = 1, N = 20, m^2 = 0.01, \lambda = 0.0017$. **Top:** We show the true Carosso RG flow of the estimated parameters of samples from ϕ^4 theory (red), as well as the same parameters computed from samples taken directly from the flow learned using our objective (5.22) (blue) and using the reverse KL objective (5.5) (teal). Our learned flow is consistently closer to the true flow than the flow learned by the reverse KL objective. The top left and center left denote the physical quantities with r_t^2 rescaling removed. Note that renormalized couplings (top right, middle right) grow with RG time, since the couplings are relevant in two dimensions. **Bottom:** We show samples along the true flow (red) and the two learned flows (blue, teal).*

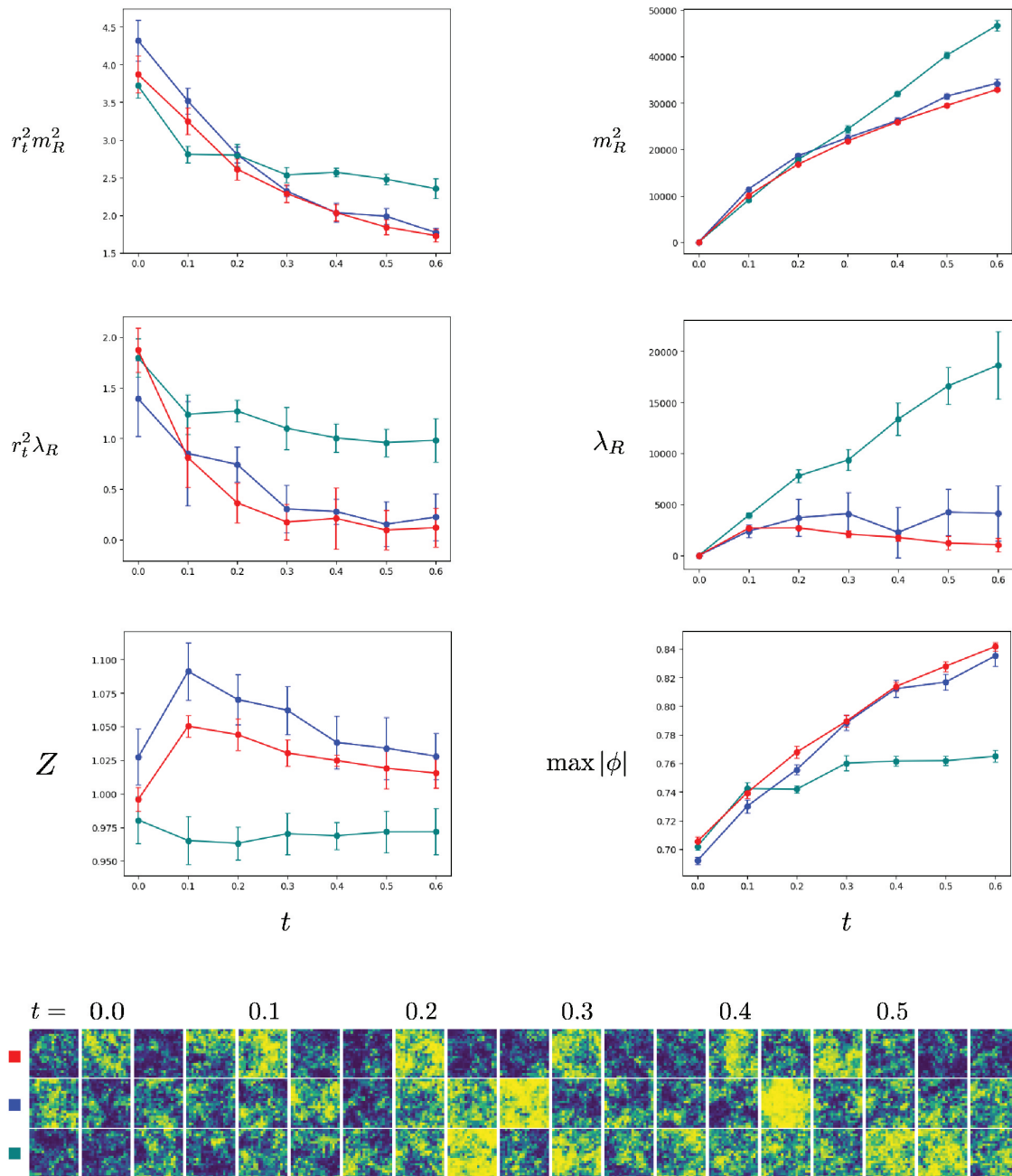


Figure 8: *Learning RG Flow of ϕ^4 theory with bare parameters $L = 1$, $N = 20$, $m^2 = -2$, $\lambda = 0.0017$. This Figure is laid out similarly to Figure 7. Note that the renormalized quartic coupling is a difficult parameter to estimate for all models, as is well-known in lattice field theory.*

8 Discussion

In this paper, we have described a methodology for designing machine-learning-based methods which *learn the RG flows of field theories*, in order to compute quantities of interest in physics. The basis of the methodology is a precise connection between the mathematics of diffusion models and exact renormalization group equations. As we have demonstrated, one can leverage this connection to design neural networks which correspond to *particular* RG schemes, have *physically-interpretable* model parameters, and yet take the form of *conventional* machine learning models. This class of models offers a powerful platform for combining physically grounded modeling intuitions with the rich numerical methods developed by the machine-learning community, and validates the hypothesis of [119] that a variational formulation of the RG flow can be used to build effective and interpretable ML models for field theory.

It is a fundamental challenge in the scientific use of ML models that ML techniques tend to lack physical interpretability, and thus are difficult to reason about and debug. This is made more challenging by a bewildering collection of numerical techniques which are comparatively difficult to evaluate systematically, and can exhibit a variety of performance characteristics which are hard to extrapolate. For example, although the performance of existing ML-based lattice samplers is promising, they indeed exhibit complex behavior [32]. Pre-existing lattice methods like the multi-grid [96, 97], cluster [98, 99], and worm [100] algorithms arise from physically-motivated heuristics. We hope that further research on equally physically-grounded ML-based sampling algorithms will lead to improved performance characteristics.

Physically-meaningful machine learning models allow for novel use cases. For example, any of the renormalizing diffusion models described in this paper generates samples from RG flowed field theory configurations during training; as such, one can use estimates of physical observables computed from batches of such configurations produced during training to diagnose how the model is learning, since the RG flows of certain physical quantities are often known from other simulations or from analytical methods. Moreover, one can conceivably modify the parameterization of the score function in a way that is interpretable in Effective Field Theory (5.26), and the coefficients of low-order terms in the Effective Field Theory thus represented by the model may also be used for diagnostics. Beyond this, the fact that the model parameters are physically meaningful should allow for hyperparameter transfer; one can think of e.g. scale setting in conventional lattice field theory as a physically-motivated method for transferring parameters between disparate models. Moreover, one should be able to use physical insight to engineer RG schemes that are particularly efficient for certain problems, including RG schemes given by nonlinear SDEs [30]. Finally, since phase transitions are *defined* by their behavior under the renormalization group, one can imagine methods which rigorously use an ML model to search for RG fixed points and RG trajectories between fixed points, thus automatically mapping out the phase diagram of a system; this would be in sharp contrast with existing ML methods for finding phase transitions, which use patterns in sampled configurations to train classifiers which can be heuristically thought of as learning order

parameters [137], rather than using the definition of a phase transition from RG.

More broadly, we see the design of machine learning methods for the physical sciences *built around the ideas of EFT* as an exciting research program. Multiscale modeling is pervasive in domains like cosmology [87] and plasma physics [138], in which effective field theories connect a series of different physical models describing a system at different scales. While our paper focuses on describing and validating the basic properties of a class of ML models which can be applied to statistical field theories, we believe that the mathematical ideas of this work should be adapted to other domains, such as hydrodynamics and molecular simulation. Indeed, the problem of rigorously learning correction terms to coarse-grained models (e.g. Maxwell-Vlasov corrections to MHD) is central to such simulations. We did not tackle those domains in this work because the formalism of the renormalization group is less fully developed in those domains, but we see this as an exciting next step, which would require developing the mathematical and computational formalism of the renormalization group beyond its most established setting. Such methods would be of considerable practical importance. An initial direction to develop such methods would be to take ‘limits’ or ‘quantizations’ of the methods described in this paper. This an exciting *theoretical* problem, not just a computational one. We hope that the union of Effective Field Theory and numerical techniques coming from machine learning can help provide a solution for the problem of unifying the traditional equation-based modeling techniques from theoretical physics and the powerful non-parametric methods of modern machine learning.

Acknowledgements

We would like to thank Daniel Ranard for valuable discussions, and Thomas Spencer for his interest in this work. JC is supported by a Junior Fellowship from the Harvard Society of Fellows, as well as in part by the Department of Energy under grant DE-SC0007870. SR is supported by NSF Mathematical Sciences Postdoctoral Fellowship, DMS-2202959.

A A brief review of literature on diffusion models

In this Appendix, we give a few brief pointers to the machine-learning literature on diffusion models for physicists interested in the subject. As mentioned in Section 1, the papers [23,25] independently introduced the basic ideas regarding diffusion models. A practical unified perspective is in [16], which forms the basis of most popular implementations of the methodology, and a more theoretical viewpoint is discussed in [17,55]. A compendium of helpful formulae useful for studying diffusion models can be found in [139].

From a practical perspective, the design of the neural network parameterizing the score function is very important, and should not be underestimated, as it encodes an implicit prior; many designs start with the U-Net [58]. Practitioners usually use comparatively simple SDEs (and it is found that the stochastic flow tends to be more efficient in practice than the deterministic flow [16]),

and care is taken with the choice of SDE solver [140]. Moreover, while structured multiscale methods [18, 20–22, 26] exist, many of the popular models such as Stable Diffusion [54] use methods which diffuse in the “space of weights of a trained convolutional neural network” [141], which allows for efficient dimensionality reduction in a way adapted to the distribution of natural images (it is easier to model a distribution over a lower-dimensional space). The literature has evolved rapidly, and commercial implementations use many carefully-tuned optimizations in order to optimize for disparate objectives like inference time and ‘image quality’, which are distinct from an accurate representation of the log-density. It is impossible to review the many applications of diffusion models to disparate domains like image, text, and audio synthesis, to scientific problems like super-resolution and denoising, and to the emergence of multimodal models; for a basic comprehensive review, we point the reader to [142].

We also caution the reader that in fact many possible paths from a background distribution to the unknown distribution can be used to design generative modeling schemes, some with improved numerical properties [143]; as such, it is not clear what aspects of the diffusive framework are essential to the effectiveness of these models. A theoretical perspective and a review of alternative ‘stochastic interpolants’ can be found in [144]. While these alternative methods offer intriguing connections to optimal transport, they have not yet been shown to be superior than simple diffusion models in applications.

B Lattice discretization of functional derivatives

Here we collect some useful conventions and notations for functional derivatives in continuum field theory, and their analogs in lattice field theory. Suppose we have a function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$, and a functional $\mathcal{F}[\phi]$. Then the functional derivative of \mathcal{F} is defined by

$$\frac{\delta}{\delta\phi(x)} \mathcal{F}[\phi(y)] := \lim_{\epsilon \rightarrow 0} \frac{\mathcal{F}[\phi(x) + \epsilon \delta^d(x-y)] - \mathcal{F}[\phi(x)]}{\epsilon}. \quad (\text{B.1})$$

For instance, if g is a differentiable function on \mathbb{R} , then if $\mathcal{G}[\phi] = \int dx g(\phi(x))$ we have $\frac{\delta\mathcal{G}[\phi]}{\delta\phi(y)} = g'(\phi(y))$.

Now suppose we consider a lattice version of $\mathcal{G}[\phi] = \int dx g(\phi(x))$, namely

$$\mathcal{G}_{\text{lattice}}[\phi(\mathbf{n})] = \left(\frac{L}{N}\right)^d \sum_{\mathbf{n} \in \mathbb{Z}_N^d} g(\phi(\mathbf{n})), \quad (\text{B.2})$$

where we are considering a d -dimensional latticization of the hypercube with side length L , such that the lattice sites are $\frac{L}{N} \mathbf{n}$ for $\mathbf{n} \in \mathbb{Z}_N^d$. We would like for the lattice analog of the functional derivative, when acting on $\mathcal{G}_{\text{lattice}}[\phi(\mathbf{n})]$, to output $g'(\phi(\mathbf{n}))$. In other words, we would like to obtain the integrand of the sum (B.2) (without the $(\frac{L}{N})^d$ factor which is an analog of the integration measure dx) with an ordinary derivative acting on the g . Then in our present notation (which

we have used throughout the body of the manuscript), the continuum functional derivative must become

$$\frac{\delta}{\delta\phi(x)} \longrightarrow \left(\frac{N}{L}\right)^d \frac{\partial}{\partial\phi(\mathbf{n})}. \quad (\text{B.3})$$

Then indeed, $\left(\frac{N}{L}\right)^d \frac{\partial}{\partial\phi(\mathbf{n})} \mathcal{G}_{\text{lattice}}[\phi(\mathbf{n})] = g'(\phi(\mathbf{n}))$. We have implicitly used the correspondence (B.3) in the manuscript to translate between equations in the continuum and equations on the lattice.

We have to be careful when comparing $\frac{\partial}{\partial\phi(\mathbf{n})}$ derivatives to their momentum space counterparts $\frac{\partial}{\partial\phi(\mathbf{p})}$. For instance, consider the lattice functional

$$\mathcal{F}_{\text{lattice}}[\phi] := \sum_{\mathbf{n}_1, \dots, \mathbf{n}_k \in \mathbb{Z}_N^d} \phi(\mathbf{n}_1) \cdots \phi(\mathbf{n}_k) f(\mathbf{n}_1, \dots, \mathbf{n}_k) \quad (\text{B.4})$$

which can also be expressed in momentum space variables as

$$\mathcal{F}_{\text{lattice}}[\tilde{\phi}] := \sum_{\mathbf{p}_1, \dots, \mathbf{p}_k \in \mathbb{Z}_N^d} \tilde{\phi}(\mathbf{p}_1) \cdots \tilde{\phi}(\mathbf{p}_k) \tilde{f}(\mathbf{p}_1, \dots, \mathbf{p}_k) \quad (\text{B.5})$$

where \tilde{f} is the Fourier transform of f . Then with the Fourier transform conventions in this paper, we have

$$\widetilde{\left[\frac{\partial \mathcal{F}_{\text{lattice}}[\phi]}{\partial\phi} \right]}(p) = \frac{1}{N^d} \frac{\partial \mathcal{F}_{\text{lattice}}[\tilde{\phi}]}{\partial\tilde{\phi}(p)}, \quad (\text{B.6})$$

where the big tilde on the left denotes that we are taking the Fourier transform of the entire left-hand side. We have used the above identity in various parts of the paper.

C Review of the Exact Renormalization Group

In this Appendix we provide a brief review of the Exact Renormalization Group (ERG) in the continuum, with an emphasis on the Wegner-Morris equation. Our expositions will follow along the lines of [119]. Consider a probability functional $P_\Lambda[\phi]$ for a scalar field theory, which captures an effective description of the system given that we can only probe momentum scales below Λ . An ERG flow addresses precisely how our effective description of the system changes as we change Λ . In particular, ERG equations take the form

$$-\Lambda \frac{d}{d\Lambda} P_\Lambda[\phi] = \mathcal{F} \left[P_\Lambda[\phi], \frac{\delta P_\Lambda[\phi]}{\delta\phi}, \frac{\delta^2 P_\Lambda[\phi]}{\delta\phi \delta\phi}, \dots \right]. \quad (\text{C.1})$$

Often it is convenient to instead work with the variable $t = -\log(\Lambda)$, in which case $-\Lambda \frac{d}{d\Lambda} \rightarrow \frac{d}{dt}$. Moreover, decreasing Λ corresponds to increasing t . An equation of the form (C.1) with an initial condition $P_{\Lambda_0}[\phi]$ prescribes how the probability functional is RG flowed for all smaller values of Λ , namely the flow determines $P_\Lambda[\phi]$ for all $0 \leq \Lambda \leq \Lambda_0$.

Consider the moment-generating function (i.e. the *partition function*) for a Euclidean scalar field theory in d spatial dimensions, namely

$$Z_\Lambda[J] := \int \mathcal{D}\phi e^{-\frac{1}{2} \int_{\mathbb{R}^d} \frac{dp}{(2\pi)^d} \left(\frac{1}{K_\Lambda(p^2)} \phi(p)\phi(-p)(p^2+m^2) + J(p)\phi(-p) \right) - S_{\text{int},\Lambda}[\phi]}, \quad (\text{C.2})$$

where $K_\Lambda(p^2)$ is a soft cutoff function (i.e. it equals 1 for $p^2 \lesssim \Lambda^2$ and is approximately 0 for $p^2 \gtrsim \Lambda^2$), and $S_{\text{int},\Lambda}[\phi]$ contains interaction terms. We will toggle between the position-space representation of the scalar field (denoted by $\phi(y)$) and its Fourier transform comprising the momentum space representation (denoted by $\phi(p)$) throughout this discussion. Note that the soft cutoff function ensures that correlation functions (at least perturbatively if the interaction terms are small) are regulated at high momentum. Further observe that by taking the functional derivative of $\log(Z_\Lambda[J])$ with respect to J we can compute moments of the probability distribution

$$P_\Lambda[\phi] \propto e^{-\frac{1}{2} \int_{\mathbb{R}^d} \frac{dp}{(2\pi)^d} \left(\frac{1}{K_\Lambda(p^2)} \phi(p)\phi(-p)(p^2+m^2) \right) - S_{\text{int},\Lambda}[\phi]}. \quad (\text{C.3})$$

It will be convenient for us to define the *action* of the above probability functional as

$$S_\Lambda[\phi] := \frac{1}{2} \int_{\mathbb{R}^d} \frac{dp}{(2\pi)^d} \left(\frac{1}{K_\Lambda(p^2)} \phi(p)\phi(-p)(p^2+m^2) \right) + S_{\text{int},\Lambda}[\phi]. \quad (\text{C.4})$$

Now suppose we are interested in flowing $P_\Lambda[\phi]$ to some smaller scale $\Lambda_R < \Lambda$, so that we only are interested in correlation functions with momentum scale $p^2 \leq \Lambda_R^2$. Further suppose that the source $J(p)$ in (C.2) satisfies $J(p) = 0$ for $p^2 > \Lambda_R^2 - \varepsilon$ for a small $\varepsilon > 0$. Then if we do not want the low-momentum correlation functions of $P_\Lambda[\phi]$ to change as we perform RG flow on the probability density, then we require

$$-\Lambda \frac{d}{d\Lambda} \log(Z_\Lambda[J]) = 0. \quad (\text{C.5})$$

A class of ERG flows which achieve this condition is given by the *Wegner-Morris* equation, namely

$$-\Lambda \frac{d}{d\Lambda} P_\Lambda[\phi] = \int_{\mathbb{R}^d} dx \frac{\delta}{\delta\phi(x)} (\Psi_\Lambda[\phi, x] P_\Lambda[\phi]), \quad (\text{C.6})$$

where $\Psi_\Lambda[\phi, x]$ is called the *reparameterization kernel*, which is required to have certain properties which we discuss below. We first note that (C.6) has a simple interpretation: as we change Λ infinitesimally by $\delta\Lambda$, the flow induces a field reparameterization as

$$\phi'(x) = \phi(x) + \frac{\delta\Lambda}{\Lambda} \Psi_\Lambda[\phi, x]. \quad (\text{C.7})$$

This is why the reparameterization kernel $\Psi_\Lambda[\phi, x]$ has its name. We see from (C.7) that in order for (C.5) to hold, we must have $\Psi_\Lambda[\phi, x]$ be a functional of ϕ that is supported at scales $p^2 \geq \Lambda^2$; that is, we are reparameterizing ϕ for momentum scales $p^2 \geq \Lambda^2$. A standard form of the reparameterization kernel is (see e.g. [14])

$$\Psi_\Lambda[\phi, x] = - \int_{\mathbb{R}^d} dy \frac{1}{2} \dot{C}_\Lambda(x-y) \frac{\delta \Sigma_\Lambda[\phi]}{\delta\phi(y)}, \quad (\text{C.8})$$

where here $\dot{C}_\Lambda(x-y)$ is a positive-definite kernel called the *ERG kernel* which is localized around $p^2 = \Lambda^2$ in momentum space, and

$$\Sigma_\Lambda[\phi] := S_\Lambda[\phi] - 2\widehat{S}_\Lambda[\phi]. \quad (\text{C.9})$$

Here $\widehat{S}_\Lambda[\phi]$ is a local action called the *seed action*, which is cut off for $p^2 \geq \Lambda^2$.

With the particular form of $\Psi_\Lambda[\phi, x]$ given (C.8), we can rewrite the Wegner-Morris equation (C.6) as a type of convection-diffusion equation, namely:

$$-\Lambda \frac{\partial P_\Lambda[\phi]}{\partial \Lambda} = \frac{1}{2} \int_{\mathbb{R}^d} dx dy \dot{C}_\Lambda(x-y) \left(\frac{\delta^2 P_\Lambda[\phi]}{\delta \phi(x) \delta \phi(y)} + 2 \frac{\delta}{\delta \phi(x)} \left(\frac{\delta \widehat{S}_\Lambda[\phi]}{\delta \phi(y)} P_\Lambda[\phi] \right) \right). \quad (\text{C.10})$$

A special case of the above gives the Polchinski flow [34]. In particular, if $\dot{C}_\Lambda(p^2)$ is the Fourier transform of the ERG kernel $\dot{C}_\Lambda(x-y)$, then the Polchinski flow is given by the specialization

$$\dot{C}_\Lambda(p^2) = (2\pi)^d (p^2 + m^2)^{-1} \Lambda \frac{\partial K_\Lambda(p^2)}{\partial \Lambda} \quad (\text{C.11})$$

$$\widehat{S}_\Lambda[\phi] = \frac{1}{2} \int_{\mathbb{R}^d} \frac{dp}{(2\pi)^d} \frac{1}{K_\Lambda(p^2)} \phi(p) \phi(-p) (p^2 + m^2). \quad (\text{C.12})$$

Notice that here, $\dot{C}_\Lambda(p^2)$ and $\widehat{S}_\Lambda[\phi]$ are both built out of the cutoff function $K_\Lambda(p^2)$ appearing in e.g. (C.2), (C.3), (C.4). The Polchinski flow has the feature that the free (i.e. Gaussian) theory

$$P_\Lambda^{\text{free}}[\phi] \propto e^{-\frac{1}{2} \int_{\mathbb{R}^d} \frac{dp}{(2\pi)^d} \left(\frac{1}{K_\Lambda(p^2)} \phi(p) \phi(-p) (p^2 + m^2) \right)}. \quad (\text{C.13})$$

satisfies the flow equation.

If we relax the assumptions on the form of the reparameterization kernel (C.8), then we can also write the Carosso scheme [8] in the form of the Wegner-Morris equation (C.6). However, the modifications to the form of the reparameterization kernel will make it so that (C.5) only approximately holds, but this is ok. Suppose that Λ_0 is the initial cutoff scale at which we begin the RG flow. The relaxed assumptions on $\Psi_\Lambda[\phi, x]$ amount to relaxed assumptions on $\dot{C}_\Lambda(x-y)$ and $\widehat{S}_\Lambda[\phi]$, namely:

1. *The cutoff function* $\dot{C}_\Lambda(p^2)$. There is a non-decreasing function $g(\Lambda)$ for $\Lambda \geq 0$ with $g(\Lambda_0) = \Lambda_0$ such that $\dot{C}_\Lambda(p^2)$ goes rapidly to zero for $|p| \geq \Lambda_0 g(\Lambda)$. We further require that $\dot{C}_\Lambda(p^2)$ is $O(1)$ for $|p| \leq \Lambda_0 g(\Lambda)$.
2. *The seed action* $\widehat{S}_\Lambda[\phi]$. Construct a normalized probability density $Q_\Lambda[\phi] \propto e^{-2\widehat{S}_\Lambda[\phi]}$. This probability density has the property that

$$\mathbb{E}_{\phi \sim Q_\Lambda[\phi]} [\phi(p_1) \phi(p_2) \cdots \phi(p_r)] \approx 0 \quad \text{for any } |p_i| \geq \Lambda_0 g(\Lambda), \quad (\text{C.14})$$

where $g(\Lambda)$ is the same function which controls the cutoff function. This means that correlation functions with momenta greater than $\Lambda_0 g(\Lambda)$ are suppressed.

Our discussion here precisely mirrors the one in Section 5.6.1, and in fact we have only lightly changed the language.

With these relaxed assumptions, the Carosso scheme is given by

$$\dot{C}_\Lambda(p^2) = e^{-p^2/\Lambda_0^2}, \quad \widehat{S}_\Lambda[\phi] = \frac{1}{2} \int_{\mathbb{R}^d} \frac{dp}{(2\pi)^d} e^{p^2/\Lambda_0^2} \phi(p)\phi(-p) p^2, \quad (\text{C.15})$$

which has $g(\Lambda) = \Lambda_0$. As explained in the body of the text, and in Carosso's paper [8], these choices allow for a simple SDE formulation of the RG flow (although again at the cost of (C.5) being only approximately satisfied).

References

- [1] F. Noé, S. Olsson, J. Köhler, and H. Wu, *Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning*, *Science* **365** (2019), no. 6457 eaaw1147, [<https://www.science.org/doi/pdf/10.1126/science.aaw1147>].
- [2] M. S. Albergo, G. Kanwar, and P. E. Shanahan, *Flow-based generative models for Markov chain Monte Carlo in lattice field theory*, 1904.1207.
- [3] R. Abbott, M. S. Albergo, A. Botev, D. Boyda, K. Cranmer, D. C. Hackett, G. Kanwar, A. G. D. G. Matthews, S. Racanière, A. Razavi, D. J. Rezende, F. Romero-López, P. E. Shanahan, and J. M. Urban, *Sampling QCD field configurations with gauge-equivariant flow models*, 2022.
- [4] M. Gerdes, P. de Haan, C. Rainone, R. Bondesan, and M. C. N. Cheng, *Learning lattice quantum field theories with equivariant continuous flows*, 2022.
- [5] D. Pfau, J. S. Spencer, A. G. D. G. Matthews, and W. M. C. Foulkes, *Ab initio solution of the many-electron Schrödinger equation with deep neural networks*, *Physical Review Res.* **2** (Sep, 2020) 033429.
- [6] J. Zinn-Justin, *Phase Transitions and Renormalization Group*. Oxford Graduate Texts. Oxford University Press, London, England, Jan., 2013.
- [7] J. L. Cardy, *Cambridge lecture notes in physics: Scaling and renormalization in statistical physics series number 5*. Cambridge University Press, Cambridge, England, Feb., 2015.
- [8] A. Carosso, *Stochastic renormalization group and gradient flow*, *Journal of High Energy Physics* **2020** (Jan., 2020).
- [9] U. Wolff, *Critical slowing down*, *Nuclear Physics B - Proceedings Supplements* **17** (1990) 93–102.

- [10] R. H. Swendsen and J.-S. Wang, *Replica Monte Carlo Simulation of Spin-Glasses*, *Physical Review Letters* **57** (Nov., 1986) 2607–2609.
- [11] A. Gelman and X.-L. Meng, *Simulating normalizing constants: from importance sampling to bridge sampling to path sampling*, *Statistical Science* **13** (May, 1998).
- [12] R. M. Neal, *Sampling from multimodal distributions using tempered transitions*, *Statistics and Computing* **6** (Dec., 1996) 353–366.
- [13] M. E. Fisher, *Renormalization group theory: Its basis and formulation in statistical physics*, *Reviews of Modern Physics* **70** (1998), no. 2 653.
- [14] O. J. Rosten, *Fundamentals of the exact renormalization group*, *Physics Reports* **511** (2012), no. 4 177–272.
- [15] W. Detmold and M. G. Endres, *Multiscale Monte Carlo equilibration: Two-color QCD with two fermion flavors*, *Physical Review D* **94** (Dec, 2016) 114502.
- [16] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models.”
- [17] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, *Score-based generative modeling through stochastic differential equations*, *arXiv:2011.13456* (2020).
- [18] E. Hoogeboom and T. Salimans, *Blurring diffusion models*, 2022.
- [19] B. Jing, G. Corso, R. Berlinghieri, and T. Jaakkola, *Subspace diffusion generative models*, 2205.0149.
- [20] F. Guth, S. Coste, V. D. Bortoli, and S. Mallat, *Wavelet score-based generative modeling*, 2022.
- [21] J. Ho, C. Saharia, W. Chan, D. J. Fleet, M. Norouzi, and T. Salimans, *Cascaded diffusion models for high fidelity image generation*, 2021.
- [22] C. Saharia, J. Ho, W. Chan, T. Salimans, D. J. Fleet, and M. Norouzi, *Image super-resolution via iterative refinement*, 2021.
- [23] J. Sohl-Dickstein, E. A. Weiss, N. Maheswaranathan, and S. Ganguli, “Deep unsupervised learning using nonequilibrium thermodynamics.”
- [24] P. Vincent, *A connection between score matching and denoising autoencoders*, . Conference Name: Neural Computation.
- [25] Y. Song and S. Ermon, *Generative modeling by estimating gradients of the data distribution*, *Advances in Neural Information Processing Systems* **32** (2019).

- [26] A. Phillips, T. Seror, M. Hutchinson, V. D. Bortoli, A. Doucet, and E. Mathieu, *Spectral diffusion processes*, 2022.
- [27] M. Lüscher, *Trivializing maps, the wilson flow and the hmc algorithm*, 0907.5491.
- [28] M. F. Atiyah and R. Bott, *The yang-mills equations over riemann surfaces*, *Philosophical Transactions of the Royal Society of London. Series A, Mathematical and Physical Sciences* **308** (1983), no. 1505 523–615.
- [29] V. G. Bornyakov, R. Horsley, R. Hudspith, Y. Nakamura, H. Perlt, D. Pleiter, P. E. L. Rakow, G. Schierholz, A. Schiller, H. Stüben, and J. M. Zanotti, *Wilson flow and scale setting from lattice qcd*, 2015.
- [30] A. Carosso, *Novel approaches to renormalization group transformations in the continuum and on the lattice*, 2020.
- [31] G. Kanwar, M. S. Albergo, D. Boyda, K. Cranmer, D. C. Hackett, S. Racanière, D. J. Rezende, and P. E. Shanahan, *Equivariant flow-based sampling for lattice gauge theory*, 2003.0641.
- [32] R. Abbott, M. S. Albergo, A. Botev, D. Boyda, K. Cranmer, D. C. Hackett, A. G. D. G. Matthews, S. Racanière, A. Razavi, D. J. Rezende, F. Romero-López, P. E. Shanahan, and J. M. Urban, *Aspects of scaling and scalability for flow-based sampling of lattice qcd*, 2022.
- [33] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, *Machine learning and the physical sciences*, *Reviews of Modern Physics* **91** (2019), no. 4 045002.
- [34] J. Polchinski, *Renormalization and effective lagrangians*, *Nuclear Physics B* **231** (1984), no. 2 269–295.
- [35] R. P. Feynman, *Difficulties in applying the variational principle to quantum field theories*, in *Variational Calculations in Quantum Field Theory: Proceedings of the International Workshop*, pp. 28–40, 1988.
- [36] S. R. White, *Density matrix formulation for quantum renormalization groups*, *Physical Review Letters* **69** (1992), no. 19 2863.
- [37] D. Perez-Garcia, F. Verstraete, M. M. Wolf, and J. I. Cirac, *Matrix product state representations*, *quant-ph/0608197* (2006).
- [38] G. Vidal, *Class of quantum many-body states that can be efficiently simulated*, *Physical Review Letters* **101** (2008), no. 11 110501.
- [39] R. N. Pfeifer, G. Evenbly, and G. Vidal, *Entanglement renormalization, scale invariance, and quantum criticality*, *Physical Review A* **79** (2009), no. 4 040301.

- [40] A. DasGupta, *Asymptotic theory of statistics and probability*, vol. 180. Springer, 2008.
- [41] T. M. Cover, *Elements of information theory*. John Wiley & Sons, 1999.
- [42] T. A. Courtade, *Monotonicity of entropy and fisher information: a quick proof via maximal correlation*, *arXiv:1610.04174* (2016).
- [43] K. P. Murphy, *Machine learning: a probabilistic perspective*. MIT press, 2012.
- [44] A. Mnih and D. Rezende, *Variational inference for Monte Carlo objectives*, in *International Conference on Machine Learning*, pp. 2188–2196, PMLR, 2016.
- [45] R. M. Neal *et. al.*, *MCMC using Hamiltonian dynamics*, *Handbook of Markov Chain Monte Carlo* **2** (2011), no. 11 2.
- [46] T. Chen, E. Fox, and C. Guestrin, *Stochastic gradient Hamiltonian Monte Carlo*, in *International Conference on Machine Learning*, pp. 1683–1691, PMLR, 2014.
- [47] M. Girolami and B. Calderhead, *Riemann manifold Langevin and Hamiltonian Monte Carlo methods*, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **73** (2011), no. 2 123–214.
- [48] M. Raginsky, A. Rakhlin, and M. Telgarsky, *Non-convex learning via stochastic gradient Langevin dynamics: a nonasymptotic analysis*, in *Conference on Learning Theory*, pp. 1674–1703, PMLR, 2017.
- [49] X. Cheng, N. S. Chatterji, P. L. Bartlett, and M. I. Jordan, *Underdamped langevin mcmc: A non-asymptotic analysis*, in *Conference on learning theory*, pp. 300–323, PMLR, 2018.
- [50] S. Chen, S. Chewi, J. Li, Y. Li, A. Salim, and A. R. Zhang, *Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions*, *arXiv:2209.11215* (2022).
- [51] E. Marinari and G. Parisi, *Simulated tempering: a new Monte Carlo scheme*, *Europhysics letters* **19** (1992), no. 6 451.
- [52] A. Gelman and X.-L. Meng, *Simulating normalizing constants: From importance sampling to bridge sampling to path sampling*, *Statistical science* (1998) 163–185.
- [53] Y. Sugita and Y. Okamoto, *Replica-exchange molecular dynamics method for protein folding*, *Chemical Physics Letters* **314** (1999), no. 1-2 141–151.
- [54] S. AI, “Stable diffusion 2.1.” <https://huggingface.co/spaces/stabilityai/stable-diffusion>, 2022. [Online; accessed Dec-2022].

- [55] Y. Song, C. Durkan, I. Murray, and S. Ermon, *Maximum likelihood training of score-based diffusion models*, *Advances in Neural Information Processing Systems* **34** (2021) 1415–1428.
- [56] Y. Song, “Generative modeling by estimating gradients of the data distribution.” <https://yang-song.net/blog/2021/score/>, 2021. [Online; accessed Dec-2022].
- [57] A. Hyvärinen, *Estimation of non-normalized statistical models by score matching*, *Journal of Machine Learning Research* **6** (2005) 695–709.
- [58] O. Ronneberger, P. Fischer, and T. Brox, *U-net: Convolutional networks for biomedical image segmentation*, in *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*, pp. 234–241, Springer, 2015.
- [59] B. D. Anderson, *Reverse-time diffusion equation models*, *Stochastic Processes and their Applications* **12** (1982), no. 3 313–326.
- [60] J. Benton, Y. Shi, V. De Bortoli, G. Deligiannidis, and A. Doucet, *From Denoising Diffusions to Denoising Markov Models*, *arXiv:2211.03595* (2022).
- [61] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud, *Neural ordinary differential equations*, *Advances in Neural Information Processing Systems* **31** (2018).
- [62] W. Grathwohl, R. T. Chen, J. Bettencourt, I. Sutskever, and D. Duvenaud, *Fjord: Free-form continuous dynamics for scalable reversible generative models*, *arXiv:1810.01367* (2018).
- [63] Z. Fodor and C. Hoelbling, *Light hadron masses from lattice QCD*, *Reviews of Modern Physics* **84** (2012), no. 2 449.
- [64] R. Bauerschmidt and T. Bodineau, *Log-Sobolev Inequality for the Continuum Sine-Gordon Model*, *Communications on Pure and Applied Mathematics* **74** (2021), no. 10 2064–2113.
- [65] J. Fröhlich, B. Simon, and T. Spencer, *Infrared bounds, phase transitions and continuous symmetry breaking*, *Communications in Mathematical Physics* **50** (1976) 79–95.
- [66] J. Glimm and A. Jaffe, *Collected Papers Vol. 1: Quantum Field Theory and Statistical Mechanics: Expositions*. Springer Science & Business Media, 1985.
- [67] D. Stroock, *Probability Theory: An Analytic View*. Cambridge University Press, 2010.
- [68] M. E. Peskin and D. V. Schroeder, *An Introduction to Quantum Field Theory*. Addison-Wesley, Reading, USA, 1995.
- [69] M. D. Schwartz, *Quantum field theory and the standard model*. Cambridge University Press, 2014.

- [70] J. V. Sengers and J. G. Shanks, *Experimental critical-exponent values for fluids*, *Journal of Statistical Physics* **137** (Oct., 2009) 857–877.
- [71] J. Cardy, *Scaling and renormalization in statistical physics*, vol. 5. Cambridge University Press, 1996.
- [72] M. Kardar, *Statistical physics of fields*. Cambridge University Press, 2007.
- [73] C. Wetterich, *Average action and the renormalization group equations*, *Nuclear Physics B* **352** (1991), no. 3 529–584.
- [74] S. Weinberg, *Phenomenological Lagrangians*, *Physica, A; (Netherlands)* **96** (1979).
- [75] H. Georgi, *Effective field theory*, *Annual Review of Nuclear and Particle Science* **43** (1993), no. 1 209–252.
- [76] A. Altland and B. D. Simons, *Condensed matter field theory*. Cambridge University Press, 2010.
- [77] T. Brauner, S. A. Hartnoll, P. Kovtun, H. Liu, M. Mezei, A. Nicolis, R. Penco, S.-H. Shao, and D. T. Son, *Snowmass white paper: effective field theories for condensed matter systems*, *arXiv:2203.10110* (2022).
- [78] S. Dubovsky, L. Hui, A. Nicolis, and D. T. Son, *Effective field theory for hydrodynamics: thermodynamics, and the derivative expansion*, *Physical Review D* **85** (2012), no. 8 085029.
- [79] M. Crossley, P. Glorioso, and H. Liu, *Effective field theory of dissipative fluids*, *Journal of High Energy Physics* **2017** (2017), no. 9 1–82.
- [80] S. Ramaswamy, *The mechanics and statistics of active matter*, *Annu. Rev. Condens. Matter Phys.* **1** (2010), no. 1 323–345.
- [81] M. C. Marchetti, J.-F. Joanny, S. Ramaswamy, T. B. Liverpool, J. Prost, M. Rao, and R. A. Simha, *Hydrodynamics of soft active matter*, *Reviews of Modern Physics* **85** (2013), no. 3 1143.
- [82] E. Weinan and B. Engquist, *Multiscale modeling and computation*, *Notices of the AMS* **50** (2003), no. 9 1062–1070.
- [83] F. Califano, G. Manfredi, and F. Valentini, *Special issue: The vlasov equation, from space to laboratory plasmas*, *Journal of Plasma Physics* **82** (2016), no. 6 701820603.
- [84] P. L. Sulem and T. Passot, *Landau fluid closures with nonlinear large-scale finite larmor radius corrections for collisionless plasmas*, *Journal of Plasma Physics* **81** (2015), no. 1 325810103.

- [85] N. I. Libeskind, R. van de Weygaert, M. Cautun, B. Falck, E. Tempel, T. Abel, M. Alpaslan, M. A. Aragón-Calvo, J. E. Forero-Romero, R. Gonzalez, S. Gottlöber, O. Hahn, W. A. Hellwing, Y. Hoffman, B. J. T. Jones, F. Kitaura, A. Knebe, S. Manti, M. Neyrinck, S. E. Nuza, N. Padilla, E. Platen, N. Ramachandra, A. Robotham, E. Saar, S. Shandarin, M. Steinmetz, R. S. Stoica, T. Sousbie, and G. Yepes, *Tracing the cosmic web*, *Monthly Notices of the Royal Astronomical Society* **473** (08, 2017) 1195–1217, [<https://academic.oup.com/mnras/article-pdf/473/1/1195/21407912/stx1976.pdf>].
- [86] C. Fidler, C. Rampf, T. Tram, R. Crittenden, K. Koyama, and D. Wands, *General relativistic corrections to N-body simulations and the Zel’dovich approximation*, *Physical Review D* **92** (Dec, 2015) 123517.
- [87] M. Vogelsberger, F. Marinacci, P. Torrey, and E. Puchwein, *Cosmological simulations of galaxy formation*, *Nature Reviews Physics* **2** (Jan., 2020) 42–66.
- [88] V. Springel, R. Pakmor, O. Zier, and M. Reinecke, *Simulating cosmic structure formation with the gadget-4 code*, *Monthly Notices of the Royal Astronomical Society* **506** (07, 2021) 2871–2949, [<https://academic.oup.com/mnras/article-pdf/506/2/2871/39271725/stab1855.pdf>].
- [89] Y. Liu, J. N. Kutz, and S. L. Brunton, *Hierarchical deep learning of multiscale differential equation time-steppers*, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* **380** (June, 2022).
- [90] S. He, Y. Li, Y. Feng, S. Ho, S. Ravanbakhsh, W. Chen, and B. Póczos, *Learning to predict the cosmological structure formation*, *Proceedings of the National Academy of Sciences* **116** (June, 2019) 13825–13832.
- [91] C. Grojean, A. Paul, Z. Qian, and I. Strümke, *Lessons on interpretable machine learning from particle physics*, *Nature Reviews Physics* **4** (2022), no. 5 284–286.
- [92] P. Lemos, N. Jeffrey, M. Cranmer, S. Ho, and P. Battaglia, *Rediscovering orbital mechanics with machine learning*, *arXiv:2202.02306* (2022).
- [93] D. Hansen, D. M. Robinson, S. Alizadeh, G. Gupta, and M. Mahoney, *Learning physical models that can respect conservation laws*, in *ICML 2023*, 2023.
- [94] G. t Hooft, *Symmetry breaking through bell-jackiw anomalies*, *Physical Review Letters* **37** (1976), no. 1 8–11.
- [95] C. G. Callan Jr, R. Dashen, and D. J. Gross, *The structure of the gauge theory vacuum*, *Physics Letters B* **63** (1976), no. 3 334–340.
- [96] J. Goodman and A. D. Sokal, *Multigrid monte carlo method. conceptual foundations*, *Physical Review D* **40** (Sep, 1989) 2035–2071.

- [97] M. G. Endres, R. C. Brower, W. Detmold, K. Orginos, and A. V. Pochinsky, *Multiscale Monte Carlo equilibration: Pure Yang-Mills theory*, *Physical Review D* **92** (2015), no. 11 114516.
- [98] R. H. Swendsen and J.-S. Wang, *Nonuniversal critical dynamics in Monte Carlo simulations*, *Physical Review Letters* **58** (1987), no. 2 86.
- [99] U. Wolff, *Collective Monte Carlo updating for spin systems*, *Physical Review Letters* **62** (1989), no. 4 361.
- [100] N. Prokof'Ev, B. Svistunov, and I. Tupitsyn, *Exact, complete, and universal continuous-time worldline Monte Carlo approach to the statistics of discrete quantum systems*, *Journal of Experimental and Theoretical Physics* **87** (1998) 310–321.
- [101] L. Dinh, J. Sohl-Dickstein, and S. Bengio, *Density estimation using real NVP*, in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, OpenReview.net, 2017.
- [102] M. S. Albergo, D. Boyda, D. C. Hackett, G. Kanwar, K. Cranmer, S. Racanière, D. J. Rezende, and P. E. Shanahan, *Introduction to normalizing flows for lattice field theory*, 2021.
- [103] J. Cotler and S. Rezchikov, *Score-based EFT, work in progress*.
- [104] D. L. Ruderman, *Origins of scaling in natural images*, *Vision research* **37** (1997), no. 23 3385–3398.
- [105] P. Mehta and D. J. Schwab, *An exact mapping between the variational renormalization group and deep learning*, *arXiv:1410.3831* (2014).
- [106] C. Bény, *Deep learning and the renormalization group*, 2013.
- [107] H. W. Lin, M. Tegmark, and D. Rolnick, *Why does deep and cheap learning work so well?*, 1608.0822.
- [108] M. Koch-Janusz and Z. Ringel, *Mutual information, neural networks and the renormalization group*, 1704.0627.
- [109] S. S. Funai and D. Giataganas, *Thermodynamics and feature extraction by machine learning*, 1810.0817.
- [110] S. Iso, S. Shiba, and S. Yokoo, *Scale-invariant feature extraction of neural network and renormalization group flow*, 1801.0717.
- [111] E. de Mello Koch, R. de Mello Koch, and L. Cheng, *Is deep learning a renormalization group flow?*, 1906.0521.

- [112] D. Albandea, L. D. Debbio, P. Hernández, R. Kenway, J. M. Rossney, and A. Ramos, *Learning trivializing flows*, 2022.
- [113] M. Lüscher, *Properties and uses of the Wilson flow in lattice QCD*, *Journal of High Energy Physics* **2010** (2010), no. 8 1–18.
- [114] S. Borsanyi, S. Dürr, Z. Fodor, C. Hoelbling, S. D. Katz, S. Krieg, T. Kurth, L. Lellouch, T. Lippert, C. McNeile, *et. al.*, *High-precision scale setting in lattice qcd*, *Journal of High Energy Physics* **2012** (2012), no. 9 1–15.
- [115] F. Wegner, *Some invariance properties of the renormalization group*, *Journal of Physics C: Solid State Physics* **7** (1974), no. 12 2098.
- [116] T. R. Morris, *A Gauge invariant exact renormalization group. 1.*, *Nucl. Phys. B* **573** (2000) 97–126, [[hep-th/9910058](#)].
- [117] J. I. Latorre and T. R. Morris, *Exact scheme independence*, *JHEP* **11** (2000) 004, [[hep-th/0008123](#)].
- [118] T. R. Morris and O. J. Rosten, *A Manifestly gauge invariant, continuum calculation of the $SU(N)$ Yang-Mills two-loop beta function*, *Physical Review D* **73** (2006) 065003, [[hep-th/0508026](#)].
- [119] J. Cotler and S. Rezchikov, *Renormalization group flow as optimal transport*, *Physical Review D* **108** (2023), no. 2 025003, [[2202.11737](#)].
- [120] D. S. Berman, M. S. Klinger, and A. G. Stapleton, *Bayesian renormalization*, *arXiv:2305.10491* (2023).
- [121] A. El Alaoui and A. Montanari, *An information-theoretic view of stochastic localization*, *IEEE Transactions on Information Theory* **68** (2022), no. 11 7423–7426.
- [122] R. Bauerschmidt, T. Bodineau, and B. Dagallier, *Stochastic dynamics and the Polchinski equation: an introduction*, *arXiv:2307.07619* (2023).
- [123] K. G. Wilson and J. B. Kogut, *The Renormalization group and the epsilon expansion*, *Phys. Rept.* **12** (1974) 75–199.
- [124] K. Fujikawa, *The gradient flow in $\lambda\phi^4$ theory*, *JHEP* **03** (2016) 021, [[1601.01578](#)].
- [125] P. Hohenberg and W. Kohn, *Inhomogeneous electron gas*, *Physical Review* **136** (1964), no. 3B B864.
- [126] W. Kohn and L. J. Sham, *Self-consistent equations including exchange and correlation effects*, *Physical Review* **140** (1965), no. 4A A1133.

- [127] W. Kohn, A. D. Becke, and R. G. Parr, *Density functional theory of electronic structure*, *The journal of physical chemistry* **100** (1996), no. 31 12974–12980.
- [128] R. Orús, *A practical introduction to tensor networks: Matrix product states and projected entangled pair states*, *Annals of physics* **349** (2014) 117–158.
- [129] N. Schuch, M. M. Wolf, F. Verstraete, and J. I. Cirac, *Computational complexity of projected entangled pair states*, *Physical Review Letters* **98** (2007), no. 14 140506.
- [130] D. P. Kingma and J. Ba, *Adam: A method for stochastic optimization*, *arXiv preprint arXiv:1412.6980* (2014).
- [131] I. Babuschkin, K. Baumli, A. Bell, S. Bhupatiraju, J. Bruce, P. Buchlovsky, D. Budden, T. Cai, A. Clark, I. Danihelka, A. Dedieu, C. Fantacci, J. Godwin, C. Jones, R. Hemsley, T. Hennigan, M. Hessel, S. Hou, S. Kapturowski, T. Keck, I. Kemaev, M. King, M. Kunesch, L. Martens, H. Merzic, V. Mikulik, T. Norman, G. Papamakarios, J. Quan, R. Ring, F. Ruiz, A. Sanchez, L. Sartran, R. Schneider, E. Sezener, S. Spencer, S. Srinivasan, M. Stanojević, W. Stokowiec, L. Wang, G. Zhou, and F. Viola, *The DeepMind JAX Ecosystem*, 2020.
- [132] U. Wolff, *Triviality of four dimensional ϕ^4 theory on the lattice*, *Scholarpedia* **9** (2014), no. 10 7367.
- [133] I. Vierhaus, *Simulation of ϕ^4 theory in the strong coupling expansion beyond the ising limit*, Master’s thesis, Humboldt-Universität zu Berlin, Mathematisch-Naturwissenschaftliche Fakultät I, 2010.
- [134] M. D. Hoffman, A. Gelman, *et. al.*, *The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo*, *J. Mach. Learn. Res.* **15** (2014), no. 1 1593–1623.
- [135] J. Lao and R. Louf, *Blackjax: A sampling library for JAX*, 2020.
- [136] Stan Development Team, *Stan modeling language users guide and reference manual*, 2023.
- [137] J. Carrasquilla and R. G. Melko, *Machine learning phases of matter*, *Nature Physics* **13** (2017), no. 5 431–434.
- [138] G. Colonna and A. D’Angola, eds., *Plasma Modeling (Second Edition)*. 2053-2563. IOP Publishing, 2022.
- [139] C. Luo, *Understanding diffusion models: A unified perspective*, 2022.
- [140] T. Karras, M. Aittala, T. Aila, and S. Laine, *Elucidating the design space of diffusion-based generative models*, *arXiv:2206.00364* (2022).

- [141] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, *High-resolution image synthesis with latent diffusion models*, in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- [142] L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, W. Zhang, B. Cui, and M.-H. Yang, *Diffusion models: A comprehensive survey of methods and applications*, 2022.
- [143] X. Liu, C. Gong, and Q. Liu, *Flow straight and fast: Learning to generate and transfer data with rectified flow*, *arXiv:2209.03003* (2022).
- [144] M. S. Albergo, N. M. Boffi, and E. Vanden-Eijnden, *Stochastic interpolants: A unifying framework for flows and diffusions*, 2023.