

GENERALIZATION THROUGH VARIANCE: HOW NOISE SHAPES INDUCTIVE BIASES IN DIFFUSION MODELS

John J. Vastola

Department of Neurobiology
Harvard Medical School
Boston, MA 02115, USA
john.vastola@hms.harvard.edu

ABSTRACT

How diffusion models generalize beyond their training set is not known, and is somewhat mysterious given two facts: the optimum of the denoising score matching (DSM) objective usually used to train diffusion models is the score function of the training distribution; and the networks usually used to learn the score function are expressive enough to learn this score to high accuracy. We claim that a certain feature of the DSM objective—the fact that its target is not the training distribution’s score, but a noisy quantity only equal to it in expectation—strongly impacts whether and to what extent diffusion models generalize. In this paper, we develop a mathematical theory that partly explains this ‘generalization through variance’ phenomenon. Our theoretical analysis exploits a physics-inspired path integral approach to compute the distributions typically learned by a few paradigmatic under- and overparameterized diffusion models. We find that the distributions diffusion models effectively learn to sample from resemble their training distributions, but with ‘gaps’ filled in, and that this inductive bias is due to the covariance structure of the noisy target used during training. We also characterize how this inductive bias interacts with feature-related inductive biases.

1 INTRODUCTION

Diffusion models (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020; Yang et al., 2023) have proven effective at producing high-quality samples (e.g., images) *like* those from some training distribution, but not overwhelmingly so. This ability to generalize is somewhat surprising for two reasons. First, the optimum of the denoising score matching (DSM) objective usually used to train them is the score function of the training distribution (Vincent, 2011; Song & Ermon, 2019), and sampling using this score only reproduces training examples (see Appendix A). Second, the network architectures usually used for score function approximation are highly expressive. Two near-SOTA models developed by Karras et al. (2022) have ~ 56 million (CIFAR-10, trained on 200 million samples) and ~ 296 million parameters (ImageNet-64, trained on 2500 million samples), respectively. Sufficiently expressive models can fit even random noise (Zhang et al., 2017).

A body of empirical work bears on the question of when and to what extent diffusion models generalize. Training data is more likely to be memorized when training sets are small (Somepalli et al., 2023a; Stein et al., 2023; Dar et al., 2024; Kadkhodaie et al., 2024), contain duplicates (Somepalli et al., 2023a; Carlini et al., 2023; Somepalli et al., 2023b), or feature low ‘complexity’ (Somepalli et al., 2023b; Stein et al., 2023). The specific training examples more likely to be memorized are either highly duplicated or outliers (Carlini et al., 2023). Whether generalization happens also strongly depends on model capacity, with Yoon et al. (2023) and Zhang et al. (2024) observing a sharp transition from memorization to generalization as the number of training examples used somewhat outstrips model capacity. However, the relationship between model performance (e.g., FID score) and model size, given a fixed number of training examples, is not monotonic; Karras et al. (2024) observe that their ImageNet models strictly improve (and hence generalize better) as model size increases.

At present, there is arguably no theory that describes when diffusion models generalize and characterizes how the associated inductive biases depend on details like training set structure, the choice

of forward/reverse processes, and model architecture. Most existing theoretical work focuses on orthogonal questions: given a *known ground truth*, can one mathematically guarantee that in some limit (e.g., a large or infinite number of samples from the ground truth distribution) diffusion models recover the ground truth, and bound how score approximation error impacts agreement (Bortoli, 2022; Chen et al., 2023a;c; Han et al., 2024)? The question we are interested in is qualitatively different: given $M \geq 1$ examples from a data distribution p_{data} , how do samples from a model trained on those examples differ from them? For example, does the model effectively interpolate training data? If so, when, and what details does this depend on? Concurrent work (Kamb & Ganguli, 2024; Niedoba et al., 2025) addresses these questions at the level of phenomenology, but not mechanism.

In this paper, we argue that six factors substantially impact how diffusion models generalize.

1. **Noisy objective.** The target of the DSM objective is not the score of the training distribution, but a noisy quantity *only equal to it in expectation*. This quantity, which we call the ‘proxy score’, introduces additional randomness to training, and has extremely high variance at low noise levels (infinite variance, in fact, at zero noise). Intuitively, this makes score function estimates, especially at low noise levels, inaccurate (this is well-known; Karras et al. (2022) remark on this when they discuss their choice of loss weighting). Moreover, this variance is not uniform in state space, but higher in ‘boundary regions’, e.g., regions of state space close to multiple training examples. This provides a useful inductive bias.
2. **Forward process.** Details of the forward process (e.g., when noise is added, asymmetry in how noise is added along different directions of state space) affect generalization through their influence on the covariance structure of the proxy score.
3. **Nonlinear score-dependence.** The learned distribution depends nonlinearly on the learned score function through the dynamics of the reverse process. This implies that the average learned distribution is *not* the training distribution, even if the score estimator is unbiased.
4. **Model capacity.** Models generalize better when # training samples \sim # model parameters.
5. **Model features.** Feature-related inductive biases interact with, and can enhance, inductive biases due to the covariance structure of the proxy score.
6. **Training set structure.** Nontrivial generalization (e.g., interpolation) is substantially more likely when a large number of training examples are near each other in state space; outliers are less likely to be meaningfully generalized.

Hence, details of training (1, 2), sampling (3), model architecture (4, 5), and the training set (6) all interact to determine the details of generalization. Other aspects, like learning dynamics, also almost certainly play a role, but we mostly neglect them here. The first factor is particularly important, and without it we will see that diffusion models do not generalize well; for this reason, we refer to the phenomenon enabled by (1) and affected by (2-6) as **generalization through variance**.

We support this claim using physics-inspired theory. The Martin-Siggia-Rose (MSR) path integral description of stochastic dynamics (Martin et al., 1973), which has also been exploited to characterize random neural networks (Crisanti & Sompolinsky, 2018) and learning dynamics (Mignacco et al., 2020; Bordelon & Pehlevan, 2022; 2023), plays a pivotal role in our analysis. First, we use the MSR path integral to derive the generic form of ‘generalization through variance’, and then we discuss in specific, analytically tractable cases of interest (e.g., linear models, lazy infinite-width neural networks) how the details change and the role of each of the aforementioned factors. To keep our theoretical analysis tractable, we focus on unconditional, non-latent models.

2 PRELIMINARIES

Data distribution. Let $p_{data}(\mathbf{x}_0)$ denote a data distribution on \mathbb{R}^D . We are especially but not exclusively interested in the case that p_{data} consists of a discrete set of $1 \leq M < \infty$ examples (e.g., images), so that $p_{data}(\mathbf{x}_0) = \sum_{m=1}^M \delta(\mathbf{x}_0 - \boldsymbol{\mu}_m)/M$, where δ is the Dirac delta function.

Forward/reverse diffusion. Training a diffusion model involves learning to convert samples from a normal distribution $p_{noise}(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{S}_T)$ to samples from $p_{data}(\mathbf{x}_0)$ via processes

$$\dot{\mathbf{x}}_t = -\beta_t \mathbf{x}_t + \mathbf{G}_t \boldsymbol{\eta}_t \quad t = 0 \rightarrow t = T \quad \text{forward process, } p_{data} \text{ to } p_{noise} \quad (1)$$

$$\dot{\mathbf{x}}_t = -\beta_t \mathbf{x}_t - \mathbf{D}_t \mathbf{s}(\mathbf{x}_t, t) \quad t = T \rightarrow t = \epsilon \quad \text{reverse process, } p_{noise} \text{ to } p_{data} \quad (2)$$

Table 1: Popular forward processes in our parameterization. For these, $\mathbf{G}_t := g_t \mathbf{I}_D$ and $\mathbf{S}_t = \sigma_t^2 \mathbf{I}_D$.

	β_t	g_t	α_t	σ_t	end time
VP-SDE	$\beta_{min} + \beta_d t$	$\sqrt{2\beta_t}$	$e^{-\int_0^t \beta_{t'} dt'}$	$\sqrt{1 - e^{-2\int_0^t \beta_{t'} dt'}}$	1
EDM	0	$\sqrt{2t}$	1	t	T

where $\boldsymbol{\eta}_t \in \mathbb{R}^K$ is Gaussian white noise, $\mathbf{G}_t \in \mathbb{R}^{D \times K}$ is a nonnegative matrix that controls the noise amplitude, $\mathbf{D}_t := \mathbf{G}_t \mathbf{G}_t^T / 2$ is the corresponding diffusion tensor, $\beta_t \geq 0$ controls decay to the origin, $\epsilon > 0$ is a time cutoff that helps ensure numerical stability, and $\mathbf{s}(\mathbf{x}, t) := \nabla_{\mathbf{x}} \log p(\mathbf{x}|t)$ is the score function. We allow \mathbf{G}_t to be a matrix so we can study how asymmetries affect generalization later. The forward process’ marginals are $p(\mathbf{x}|t) := \int p(\mathbf{x}|\mathbf{x}_0, t) p_{data}(\mathbf{x}_0) d\mathbf{x}_0$. The transition probabilities are $p(\mathbf{x}|\mathbf{x}_0, t) = \mathcal{N}(\mathbf{x}; \alpha_t \mathbf{x}_0, \mathbf{S}_t)$, where $\alpha_t := e^{-\int_0^t \beta_{t'} dt'}$ and $\mathbf{S}_t := \int_0^t 2\mathbf{D}_{t'} \alpha_{t'}^2 dt'$.

The forward process assumed here is fairly general, and includes popular choices like the VP-SDE (Song et al., 2021) and EDM formulation (Karras et al., 2022) (Table 1). This choice of reverse process is called the probability flow ODE (PF-ODE), and has been shown to have both practical (Song et al., 2021) and theoretical (Chen et al., 2023b) advantages. Since $\mathbf{s}(\mathbf{x}, t)$ is required to run the reverse process but is a priori unknown, ‘‘training’’ a model means approximating $\mathbf{s}(\mathbf{x}, t)$.

Denosing score matching. Given $P \gg 1$ independent samples from $p(\mathbf{x}, \mathbf{x}_0, t)$ (note: P is different than M , the number of points in discrete p_{data}), one could use a mean-squared-error objective

$$J_0(\boldsymbol{\theta}) := \mathbb{E}_{t, \mathbf{x}} \left\{ \frac{\lambda_t}{2} \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{s}(\mathbf{x}, t)\|_2^2 \right\} = \int \frac{\lambda_t}{2} \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{s}(\mathbf{x}, t)\|_2^2 p(\mathbf{x}|t) p(t) d\mathbf{x} dt \quad (3)$$

to learn a parameterized score estimator $\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t)$. Here, $\lambda_t > 0$ is a positive weighting function and $p(t)$ is a time-sampling distribution. The DSM objective (Vincent, 2011; Song & Ermon, 2019)

$$J_1(\boldsymbol{\theta}) := \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}} \left\{ \frac{\lambda_t}{2} \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)\|_2^2 \right\} = \int \frac{\lambda_t}{2} \|\hat{\mathbf{s}}_{\boldsymbol{\theta}} - \tilde{\mathbf{s}}\|_2^2 p(\mathbf{x}, \mathbf{x}_0, t) d\mathbf{x} d\mathbf{x}_0 dt \quad (4)$$

where $p(\mathbf{x}, \mathbf{x}_0, t) := p(\mathbf{x}|\mathbf{x}_0, t) p_{data}(\mathbf{x}_0) p(t)$, is usually used instead. While the folklore justifying this choice is that the score function is not known, this is not true; both J_0 and J_1 are optimized when $\hat{\mathbf{s}}_{\boldsymbol{\theta}}$ equals the score of the training distribution (see Appendix A), which is known.

We will argue that the real difference between J_0 and J_1 is that J_1 generalizes better, and that this is in part because the **proxy score** $\tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0) := \nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{x}_0, t) = \mathbf{S}_t^{-1}(\alpha_t \mathbf{x}_0 - \mathbf{x})$ is used as the target instead of the true score. It is a ‘noisy’ version of the true score (see Appendix B), since

$$\mathbb{E}_{\mathbf{x}_0|\mathbf{x}, t}[\tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)] = \mathbf{s}(\mathbf{x}, t) \quad C_{ij}(\mathbf{x}, t) := \text{Cov}_{\mathbf{x}_0|\mathbf{x}, t}[\tilde{s}_i, \tilde{s}_j] = S_{t, ij}^{-1} + \partial_{ij}^2 \log p(\mathbf{x}|t). \quad (5)$$

Although the proxy score is equal to the score of the training distribution in expectation, neural networks trained on J_1 empirically learn a different distribution and generalize better. We claim that this fact is closely related to the *covariance structure* of the proxy score. Two relevant observations about its form are as follows. First, it is large at small times, since $\mathbf{S}_t \rightarrow \mathbf{0}$ as $t \rightarrow 0$. Second, it is large where the log-likelihood $\log p(\mathbf{x}|t)$ has substantial curvature. In the typical case, where p_{data} consists of a discrete set of M examples, regions of high curvature correspond to the location of training examples and the boundaries between them (Fig. 1; see Appendix C for more discussion).

Generalization and inductive biases. In a typical supervised learning setting, one trains a model on one set of data and tests it on another, and defines ‘generalization error’ as performance on the held-out data. Here, we are interested in a different type of problem: *given a model trained on samples from $p(\mathbf{x}, \mathbf{x}_0, t)$, to what extent does the learned distribution differ from p_{data} , and what are the associated inductive biases?* Of particular interest is whether models do three things: (i) interpolation (filling in gaps in the training data), (ii) extrapolation (extending patterns in the training data), and (iii) feature blending (generating samples which include both feature X and feature Y even when training examples only involve one of the two features).

In our setting, a subtle but important point is that there is generally no ground truth. For example, the smooth distribution that CIFAR-10 or MNIST images are drawn from does not exist, except in a ‘Platonic’ sense; we are interested in the extent to which diffusion models learn a distribution plausibly *like* a smoothed version of the training distribution.

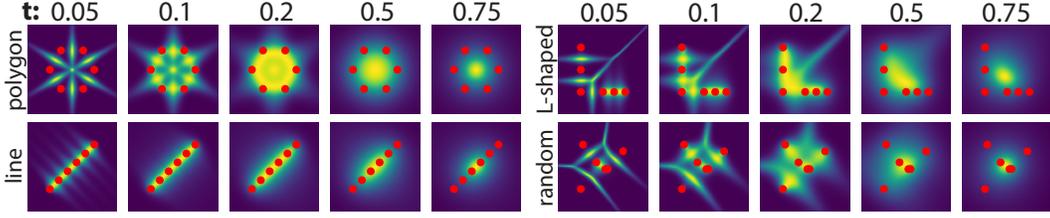


Figure 1: Visualization of proxy score variance ($\text{tr}(\mathbf{C})/[\text{tr}(\mathbf{C}) + \|\mathbf{s}\|_2^2]$) for four example 2D data distributions. Each example data distribution is supported on a small number of point masses (red dots). As t changes (left: small t , right: large t), boundary regions at different scales are emphasized.

3 APPROACH: COMPUTING TYPICAL LEARNED DISTRIBUTIONS

The distribution $q(\mathbf{x}_0|\boldsymbol{\theta})$ learned by a diffusion model depends on the learned score $\hat{\mathbf{s}}_\theta$ nonlinearly through PF-ODE dynamics; importantly, we are less interested in how well the score is estimated, and more interested in how estimation errors impact q . The learned score can be viewed as a random variable, since it depends on the P samples $\mathbf{x}^{(i)}, \mathbf{x}_0^{(i)}, t^{(i)} \sim p(\mathbf{x}, \mathbf{x}_0, t)$ used during training. In order to theoretically understand how diffusion models generalize, we aim to obtain an analytic expression for the ‘typical’ learned distribution by averaging q over sample realizations.

How do we do the required averaging? One of our major contributions is to introduce a theoretical approach for averaging $q(\mathbf{x}_0)$ over variation due to $\hat{\mathbf{s}}$. Below, we describe our approach.

Writing PF-ODE dynamics in terms of a path integral. How does one average over the result of an ODE given that, in the case of PF-ODE dynamics, there is generally no closed-form expression for the result? To address this issue, we use a novel **stochastic path integral** representation of PF-ODE dynamics that makes the required average easy to do. If $q(\mathbf{x}_0|\mathbf{x}_T; \boldsymbol{\theta})$ denotes the distribution of PF-ODE outputs given a score estimator $\hat{\mathbf{s}}_\theta(\mathbf{x}, t)$ and a fixed noise seed \mathbf{x}_T ,

$$q(\mathbf{x}_0|\mathbf{x}_T; \boldsymbol{\theta}) = \int \mathcal{D}[\mathbf{p}_t] \mathcal{D}[\mathbf{x}_t] \exp \left\{ \int_\epsilon^T i \mathbf{p}_t \cdot [\dot{\mathbf{x}}_t + \beta_t \mathbf{x}_t + \mathbf{D}_t \hat{\mathbf{s}}_\theta(\mathbf{x}_t, t)] dt \right\} \quad (6)$$

where the integral is over all possible paths from \mathbf{x}_T to \mathbf{x}_0 . (To avoid technical issues, we assume a particular time discretization in all calculations. See Appendix D.) This type of path integral is a time-reversed version of the Martin-Siggia-Rose (MSR) path integral (Martin et al., 1973).

Averaging over possible sample realizations. Because the argument of the exponential depends linearly on the score, the required ensemble average is now easy to do. Using $[\dots]$ to denote it,

$$[q(\mathbf{x}_0|\mathbf{x}_T)] = \int \mathcal{D}[\mathbf{p}_t] \mathcal{D}[\mathbf{x}_t] \exp \left\{ M_1 - \frac{1}{2} M_2 + \dots \right\} \quad (7)$$

$$M_1 := \int_\epsilon^T i \mathbf{p}_t \cdot [\dot{\mathbf{x}}_t + \beta_t \mathbf{x}_t + \mathbf{D}_t \mathbf{s}_{avg}(\mathbf{x}_t, t)] dt \quad M_2 := \int_\epsilon^T \int_\epsilon^T \mathbf{p}_t^T \mathbf{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') \mathbf{p}_{t'} dt dt'$$

where $\mathbf{s}_{avg}(\mathbf{x}_t, t) := [\hat{\mathbf{s}}_\theta(\mathbf{x}_t, t)]$ is the ensemble’s average score estimator, and $\mathbf{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') := \mathbf{D}_t \text{Cov}_\theta[\hat{\mathbf{s}}(\mathbf{x}_t, t), \hat{\mathbf{s}}(\mathbf{x}_{t'}, t')] \mathbf{D}_{t'}$ measures ensemble variance. Assuming higher-order terms can be neglected—and hence that the estimator distribution is approximately Gaussian—one can show (see Appendix D) that sampling from $[q(\mathbf{x}_0|\mathbf{x}_T)]$ is equivalent to integrating an (Ito-interpreted) SDE:

Proposition 3.1 (Effective SDE description of typical learned distribution). *Sampling from $[q(\mathbf{x}_0|\mathbf{x}_T)]$ is equivalent to integrating the (Ito-interpreted) SDE*

$$\dot{\mathbf{x}}_t = -\beta_t \mathbf{x}_t - \mathbf{D}_t \mathbf{s}_{avg}(\mathbf{x}_t, t) + \boldsymbol{\xi}(\mathbf{x}_t, t) \quad t = T \rightarrow t = \epsilon \quad (8)$$

with initial condition \mathbf{x}_T , where $\mathbf{s}_{avg}(\mathbf{x}_t, t) := [\hat{\mathbf{s}}_\theta(\mathbf{x}_t, t)]$ and where the noise term $\boldsymbol{\xi}(\mathbf{x}_t, t)$ has mean zero and autocorrelation $\mathbf{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') := \mathbf{D}_t \text{Cov}_\theta[\hat{\mathbf{s}}(\mathbf{x}_t, t), \hat{\mathbf{s}}(\mathbf{x}_{t'}, t')] \mathbf{D}_{t'}$.

If $\hat{\mathbf{s}}$ is unbiased and M is finite, then the noise term is solely responsible for the difference between true PF-ODE dynamics (which reproduces training examples) and a model’s ‘typical’ sampling

dynamics—i.e., generalization occurs if and only if $\mathbf{V} \neq \mathbf{0}$. This makes characterizing \mathbf{V} , which we call the *V-kernel* since it reflects ensemble variance, crucially important for understanding how diffusion models generalize. Our remaining theoretical work is to complete two tasks: first, to compute s_{avg} and \mathbf{V} for a few paradigmatic and theoretically tractable architectures; and second, to study how their precise forms affect $[q(\mathbf{x}_0)]$.

4 DIFFUSION MODELS THAT MEMORIZE TRAINING DATA STILL GENERALIZE

It is instructive to first consider an extreme case: do diffusion models generalize in the complete *absence* of any model-related inductive biases? Perhaps surprisingly, the answer is yes. In this section, we make this point using a toy model in which training and sampling are interleaved.

Suppose the PF-ODE is integrated backward in time from an initial point \mathbf{x}_T until $t = \epsilon$ using first-order Euler updates of size Δt . At each time step, suppose one samples $\mathbf{x}_{0t} \sim p(\mathbf{x}_0|\mathbf{x}_t, t) = \frac{p(\mathbf{x}_t|\mathbf{x}_0, t)p_{data}(\mathbf{x}_0)}{p(\mathbf{x}_t|t)}$, constructs the ‘naive’ score estimator $\hat{\mathbf{s}}(\mathbf{x}_t, t) := \mathbf{s}(\mathbf{x}_t, t) + \sqrt{\frac{\kappa}{\Delta t}}[\tilde{\mathbf{s}}(\mathbf{x}_t, t; \mathbf{x}_{0t}) - \mathbf{s}(\mathbf{x}_t, t)]$, and uses this estimator as the score for that update. Assume this process continues, with new samples drawn at each time step. Despite this approach using the proxy score directly (so that training data is ‘memorized’), one obtains a nontrivial V-kernel, and hence generalization:

Proposition 4.1 (Naive score estimator generalizes). *Consider the result of integrating the PF-ODE (Eq. 2) from $t = T$ to $t = \epsilon$ using first-order Euler updates of the form*

$$\mathbf{x}_{t-1} = \mathbf{x}_t + \Delta t \left\{ \beta_t \mathbf{x}_t + \mathbf{D}_t \left(\mathbf{s}(\mathbf{x}_t, t) + \sqrt{\frac{\kappa}{\Delta t}} [\tilde{\mathbf{s}}(\mathbf{x}_t, t; \mathbf{x}_{0t}) - \mathbf{s}(\mathbf{x}_t, t)] \right) \right\}, \quad \mathbf{x}_{0t} \sim p(\mathbf{x}_0|\mathbf{x}_t, t).$$

Then $[q(\mathbf{x}_0|\mathbf{x}_T)]$ is described by an effective SDE (Eq. 8) with $s_{avg} = \mathbf{s}$ and V-kernel

$$\mathbf{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') := \kappa \mathbf{D}_t \mathbf{C}(\mathbf{x}, t) \mathbf{D}_{t'} \delta(t - t'). \quad (9)$$

See Appendix E for details. Notably, the effective SDE is noisier when the covariance \mathbf{C} of the proxy score is high, e.g., in boundary regions between training examples. Next, we will see that this is also true for less trivial models, but that the proxy score’s covariance interacts with feature-related biases in order to determine the SDE’s overall noise term.

5 FEATURE-RELATED INDUCTIVE BIASES MODULATE GENERALIZATION

Model architecture is known to produce certain inductive biases, with spectral bias being a well-known example (Rahaman et al., 2019; Bordelon et al., 2020; Canatar et al., 2021). How do model-feature-related inductive biases affect the V-kernel? We answer this question below in two interesting but tractable cases: linear models, and (lazy regime) infinite-width neural networks.

5.1 THE V-KERNEL OF EXPRESSIVE LINEAR MODELS

In what follows, we may write $\mathbf{z} := (\mathbf{x}, t)$ to ease notation. Consider a linear score estimator

$$\hat{\mathbf{s}}_{\theta}(\mathbf{x}, t) = \mathbf{w}_0 + \mathbf{W} \phi(\mathbf{x}, t), \quad (10)$$

where the F feature maps $\phi := (\phi_1, \dots, \phi_F)^T$ are linearly independent, smooth functions from $\mathbb{R}^D \times (0, T]$ to \mathbb{R} that are square-integrable with respect to the measure $\lambda_t p(\mathbf{x}, t)$. The parameters to be estimated are $\theta := \{\mathbf{w}_0, \mathbf{W}\}$, with $\mathbf{w}_0 \in \mathbb{R}^D$ and $\mathbf{W} \in \mathbb{R}^{D \times F}$. Note that this estimator is linear in its features, but not necessarily in \mathbf{x} or t . The weights that optimize Eq. 4 are (see Appendix F)

$$\mathbf{W}^* = -\mathbf{J}^T \Sigma_{\phi}^{-1} \quad \mathbf{w}_0^* = \mathbf{J}^T \Sigma_{\phi}^{-1} \langle \phi \rangle + \langle \tilde{\mathbf{s}} \rangle \quad (11)$$

where we define $\langle \dots \rangle := \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t}[\lambda_t \dots] / \mathbb{E}_t[\lambda_t]$ and matrices

$$\mathbf{J} := -\langle [\phi(\mathbf{x}, t) - \langle \phi \rangle] [\tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0) - \langle \tilde{\mathbf{s}} \rangle]^T \rangle \quad \Sigma_{\phi} := \langle [\phi(\mathbf{x}, t) - \langle \phi \rangle] [\phi(\mathbf{x}, t) - \langle \phi \rangle]^T \rangle.$$

When averaged over \mathbf{x}_0 sample realizations, the estimator $\hat{\mathbf{s}}_*(\mathbf{x}, t) = \mathbf{w}_0^* + \mathbf{W}^* \phi(\mathbf{x}, t)$ is unbiased as long as the set of feature maps is sufficiently expressive. Interestingly, this is true regardless of the \mathbf{x} or t samples used, provided $F \leq P$. The following result characterizes $[q(\mathbf{x}_0)]$ for linear models:

Proposition 5.1 (Expressive linear models asymptotically generalize). *Suppose the parameters of an expressive linear score estimator (Eq. 10) with F features are perfectly optimized according to the DSM objective (Eq. 4) using P independent samples from $p(\mathbf{x}, \mathbf{x}_0, t)$. Define the feature kernel*

$$k(\mathbf{z}; \mathbf{z}') := \frac{1}{\sqrt{F}} [\phi(\mathbf{z}) - \langle \phi \rangle]^T \Sigma_\phi^{-1} [\phi(\mathbf{z}') - \langle \phi \rangle]. \quad (12)$$

Let $\kappa := F/P$. Provided that the limit exists and is finite, in the $P \rightarrow \infty$ limit, we have

$$\mathbf{V}(\mathbf{z}; \mathbf{z}') = \lim_{P \rightarrow \infty} \kappa \mathbf{D}_t \mathbb{E}_{\mathbf{z}''} \left\{ \frac{\lambda_{t''}^2}{\mathbb{E}_t[\lambda_t]^2} k(\mathbf{z}; \mathbf{z}'') \mathbf{C}(\mathbf{z}'') k(\mathbf{z}''; \mathbf{z}') \right\} \mathbf{D}_{t'}. \quad (13)$$

Note that if the number of features F does not scale with P , $\mathbf{V} \equiv \mathbf{0}$. See Appendix F for the details of our argument. The V-kernel for linear models differs from the naive score’s V-kernel (Eq. 9) via the presence of feature-related factors—in particular, the effective SDE is noisier *where features take atypical values*. One expects that these factors can either enhance or compete with noise due to the covariance structure (e.g., noise is higher if features take atypical values in boundary regions).

5.2 THE V-KERNEL OF LAZY INFINITE-WIDTH NEURAL NETWORKS

Neural networks in the neural tangent kernel (NTK) regime (Jacot et al., 2018; Bietti & Mairal, 2019) provide another interesting but tractable model. Such networks exhibit ‘lazy’ learning (Chizat et al., 2019) in the sense that weights do not move much from their initial values. Moreover, it is known that they interpolate training data in the absence of parameter regularization or early stopping (Bordelon et al., 2020). If they precisely interpolated their samples, we would expect to recover a V-kernel like the one we computed in Sec. 4; more generally, we expect something similar modified by the spectral inductive biases associated with the architecture (Canatar et al., 2021).

For simplicity, we consider fully-connected networks whose hidden layers all have width N , which is taken to infinity together with P (see Appendix G for details). The associated NTK has a Mercer decomposition with respect to the measure $\lambda_t p(\mathbf{x}, t) / \mathbb{E}[\lambda_t]$, so K can be written in terms of F orthonormal features $\{\phi_i\}$:

$$K(\mathbf{x}, t; \mathbf{x}', t') = \sum_k \lambda_k \phi_k(\mathbf{x}, t) \phi_k(\mathbf{x}', t') \quad \int \frac{\lambda_t}{\mathbb{E}[\lambda_t]} \phi_k(\mathbf{x}, t) \phi_\ell(\mathbf{x}, t) p(\mathbf{x}, t) d\mathbf{x} dt = \delta_{k\ell}. \quad (14)$$

We assume training involves full-batch gradient descent on P samples from $p(\mathbf{x}, \mathbf{x}_0, t)$, so that the learned score function after training for an amount of ‘time’ τ has the closed-form solution

$$\hat{\mathbf{s}}(\mathbf{z}) = \hat{\mathbf{s}}_0(\mathbf{z}) + [\tilde{\mathbf{S}} - \hat{\mathbf{S}}_0]^T (\mathbf{I} - e^{-\Lambda_\tau \mathbf{K} \tau / P}) \mathbf{K}^{-1} \mathbf{k}(\mathbf{z})$$

where $\hat{\mathbf{s}}_0$ is the network’s initial output, $\tilde{\mathbf{S}} \in \mathbb{R}^{D \times P}$ contains proxy score samples, $\hat{\mathbf{S}}_0 \in \mathbb{R}^{D \times P}$ contains the network’s initial outputs given the samples, $\mathbf{K} \in \mathbb{R}^{P \times P}$ is the kernel Gram matrix, $\Lambda_\tau \in \mathbb{R}^{P \times P}$ is a diagonal matrix containing the weighting function $\lambda_t / \mathbb{E}[\lambda_t]$ evaluated on samples, and $\mathbf{k}(\mathbf{x}, t)$ is an input-dependent vector whose i -th component is $K(\mathbf{x}^{(i)}, t^{(i)}; \mathbf{x}, t)$. We have:

Proposition 5.2 (Lazy neural networks asymptotically generalize). *Suppose the parameters of a fully-connected, infinite-width neural network characterized by a rank F NTK are optimized according to the DSM objective (Eq. 4) using P independent samples from $p(\mathbf{x}, \mathbf{x}_0, t)$ via full-batch gradient descent for a training ‘time’ τ . Define the feature kernel*

$$k(\mathbf{z}; \mathbf{z}') := \frac{1}{\sqrt{F}} \phi(\mathbf{z})^T (\mathbf{I}_F - e^{-\Lambda \tau}) \phi(\mathbf{z}'). \quad (15)$$

Let $\kappa := F/P$. Provided that the limit exists and is finite, in the $P \rightarrow \infty$ limit, we have

$$\mathbf{V}(\mathbf{z}; \mathbf{z}') = \lim_{P \rightarrow \infty} \kappa \mathbf{D}_t \mathbb{E}_{\mathbf{z}''} \left\{ \frac{\lambda_{t''}^2}{\mathbb{E}[\lambda_t]^2} k(\mathbf{z}; \mathbf{z}'') \mathbf{C}(\mathbf{z}'') k(\mathbf{z}''; \mathbf{z}') \right\} \mathbf{D}_{t'}. \quad (16)$$

In the infinite training time limit, we recover the Sec. 4 result with a prefactor $\kappa(\Delta \mathbf{z}) = \text{const.}$:

$$\mathbf{V}(\mathbf{z}; \mathbf{z}') = \kappa(\Delta \mathbf{z}) \mathbf{D}_t \mathbf{C}(\mathbf{z}) \mathbf{D}_t \delta(\mathbf{z} - \mathbf{z}'). \quad (17)$$

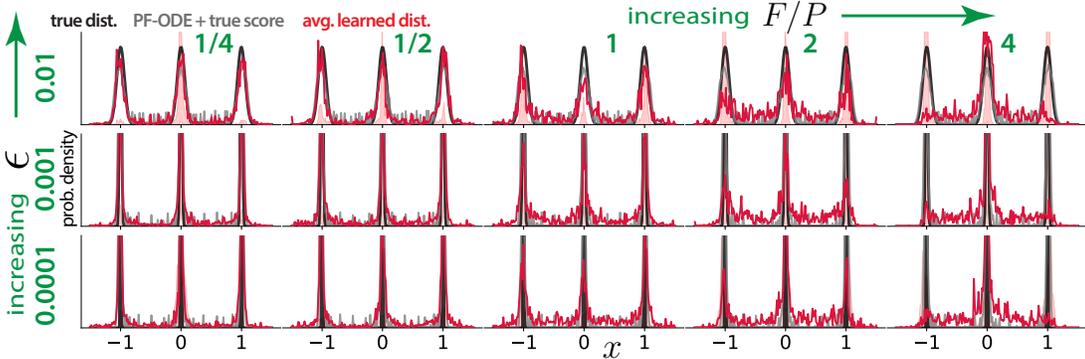


Figure 2: Average learned distribution ($N = 100$) for a linear model with Gaussian features trained on different sample draws from a 1D data distribution $\{-1, 0, 1\}$. Red: average learned distribution; black: true distribution; gray: PF-ODE approximation of true distribution. Different values of the time cutoff ϵ and ratio F/P are shown. Note that there is more generalization as both become larger.

See Appendix G for the full details of our argument. Interestingly, although the network is not assumed to be in the feature-learning regime, this result interpolates between our pure memorization (Prop. 4.1) and linear model (Prop. 5.1) results as we change the value of the training time τ . The feature-related inductive biases that appear are precisely the well-known spectral biases.

See Appendices H and I for discussion of how to obtain analogous results for diffusion models with slightly different training objectives, like those for which a ‘denoiser’ rather than a score approximator is learned (see, e.g., Karras et al. (2022)).

6 GENERALIZATION THROUGH VARIANCE: CONSEQUENCES AND EXAMPLES

In this section, we briefly discuss salient consequences of generalization through variance.

Benign properties of generalization through variance. In what sense might generalization through variance provide a ‘reasonable’ inductive bias? Its key driver is the proxy score covariance, which is large primarily in boundary regions between training examples (Appendix C), and this fact greatly constrains the way this type of generalization can occur. A data set with one data point (so that $M = 1$) is not generalized, since the proxy score covariance is trivially zero. If the data distribution is primarily supported on some low-dimensional ‘data manifold’, the proxy score covariance tends to be nontrivial only along that manifold, and hence generalization through variance preserves the dimensionality of the data manifold.

Very far from training examples, the proxy score covariance is approximately zero, so there is no generalization through variance. Finally, the effective PF-ODE both follows deterministic PF-ODE dynamics *on average*, since the V-kernel-related noise term has mean zero, and is also *most likely* to follow deterministic PF-ODE dynamics, since the probability of paths that deviate from it can be shown to be somewhat lower. This means that, although effective PF-ODE dynamics differ from the deterministic PF-ODE’s dynamics, they do not *substantially* differ, meaning that regions near training data will still tend to be sampled most. See Appendix J for more details and discussion.

Memorization and the V-kernel in the small noise limit. Our characterization of the typical learned distribution $[q]$ in terms of a stochastic process is somewhat unsatisfying, in part because it remains unclear how the V-kernel affects the way $[q]$ generalizes p_{data} . One can make some progress on the issue by making a small noise approximation, which is valid (for example) when models are somewhat underparameterized, so that $\kappa = F/P$ is somewhat less than 1. When the effective PF-ODE’s noise term is sufficiently small, one can invoke a semiclassical approximation of the relevant path integral. We find (Appendix K) that, at least in this limit,

$$[q(\mathbf{x}_0)] \approx p(\mathbf{x}_0|\epsilon) \frac{1}{\sqrt{\det\left(\frac{1}{\kappa} \frac{\partial^2 \mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0))}{\partial \mathbf{x}_T \partial \mathbf{x}_T}\right)}} \quad (18)$$

where \mathcal{S}_{cl} quantifies the (negative log-) likelihood of the most likely path that goes from a noise seed x_T to a sample x_0 . In words: $[q]$ equals the (ϵ -noise-corrupted) data distribution, times a curvature factor that quantifies the likelihood of small deviations from deterministic PF-ODE dynamics. It is through the V-kernel’s influence on this curvature term that it affects generalization (although unfortunately it appears difficult to be more explicit about how it does so, at least analytically).

Gap-filling inductive bias. Given that the V-kernel is especially sensitive to the ‘gaps’ between training examples, one expects that generalization through variance works by effectively filling in these gaps. This appears to be often, but not always, true. First, in its naive form (see, e.g., Sec. 4) generalization through variance can actually *reduce* the probability associated with boundary regions, since the additional noise in those regions makes the dynamics spend less time in them. If there is nontrivial temporal generalization, for example via time-dependent features ϕ , the V-kernel may have a nontrivial temporal autocorrelation structure; we speculate that these autocorrelations may be a key mechanism that allows the dynamics to spend more time in boundary regions.

Second, the details of generalization are strongly modulated by two numbers: the time cutoff ϵ , and the ratio F/P that determines the extent to which a model is over- or underparameterized. Fig. 2 depicts an illustrative one-dimensional example where there are three training examples $\{-1, 0, 1\}$, and where the model is linear (see Prop. 5.1) with Gaussian features centered at different values of x and t , each with the same width. The average learned distribution (red) tends to differ from both the true distribution (black) and its PF-ODE approximation (gray) in the size of peaks near training data, and in the regions between training data. These differences are larger when F/P and ϵ are larger. Taking both large produces the largest difference, but not obviously the ‘best’ generalization of training data.

Feature-noise alignment affects generalization. Different feature sets interact with the structure of the proxy score covariance differently, and hence produce different kinds of generalization. Fig. 3 shows how the same 2D data distribution (four examples, which together determine the vertices of a square) is generalized differently depending on its orientation, and depending on which linear model feature set (here, either Gaussian or Fourier features) is used.

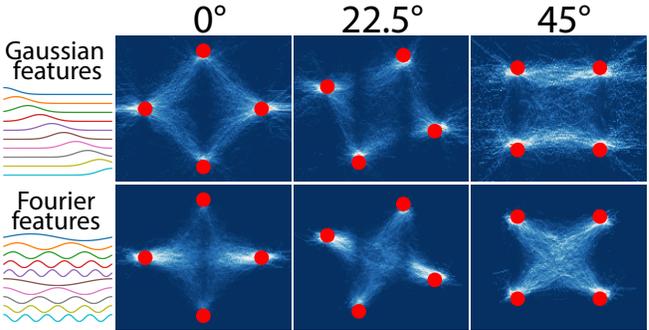


Figure 3: Generalization of a 2D data distribution depends on features used and data orientation. Heatmaps of samples from $N = 100$ linear models are shown in different conditions, with training data (red dots) overlaid. Notice that which gaps are ‘filled in’, e.g., whether a square shape or cross shape is made, depends on both factors.

7 DISCUSSION

We used a novel path-integral approach to quantitatively characterize the ‘typical’ distribution learned by diffusion models, and find that generalization is influenced by a combination of factors related to training (the DSM objective and forward process; Sec. 2 and 4), sampling (the learned distribution depends nonlinearly on score estimates; Sec. 3), model architecture (Sec. 5), and the data distribution. Below, we use our theory to comment on various previous observations.

DSM produces noisy estimators, but stable distributions. Various forms of score ‘mislearning’ are well-known. At small times, scores are hard to learn due to the noisiness of the proxy score target, leading authors like Karras et al. (2022) to suggest a $p(t)$ that emphasizes intermediate noise scales. Chao et al. (2022) discuss how score estimation errors affect conditional scores. Xu et al. (2023) explicitly study the variance-near-mode-boundaries issue we discussed, and propose a strategy for mitigating it. On the other hand, it is well-known that despite noisy score estimates, diffusion models generally produce smooth output distributions (see, e.g., Luzzi et al. (2024)). Moreover, two diffusion models trained on non-overlapping subsets of a data set are often highly similar (Kadkhodaie et al., 2024). These facts are due to noisy score estimates contributing to sample generation through the PF-ODE, which effectively ‘averages’ over estimator noise. Our theory is consistent with these

observations: even the interleaved training-sampling procedure discussed in Sec. 4 produces a well-behaved, smooth distribution.

DSM produces a boundary-smearing inductive bias. This has been previously pointed out by authors like Xu et al. (2023). Where we differ from previous authors is in considering this issue a potential strength. Integrating the PF-ODE using the true score reproduces training examples, so it is in some sense beneficial to ‘mislearn’ the score. This particular kind of mislearning is useful for several ways of generalizing point clouds, including interpolation, extrapolation, and feature blending. Moreover, producing this inductive bias is an interesting way diffusion models differ from something like kernel density estimation: boundary regions *across different noise scales* are smeared out, with different scales linked via PF-ODE dynamics, which may provide better generalization than convolving the training distribution with any single kernel.

Architecture-related inductive biases play a role. As we showed in Sec. 5, feature/architecture-related inductive biases interact with DSM’s boundary-smearing bias in order to determine how diffusion models generalize. This appears to be consistent, for example, with the Kadkhodaie et al. (2024) finding that diffusion models effectively exhibit ‘adaptive geometric harmonic priors’; their finding is specifically in the context of score estimation using a convolutional neural network (CNN) architecture. It is plausible that this choice encourages a harmonic inductive bias, since CNNs more generally exhibit inductive biases related to translation equivariance (Cohen & Welling, 2016).

Generalization through variance harmful and helpful. It is important to note that this kind of generalization is not always helpful. A trivial example is that unconditional models trained on MNIST digit images tend to learn to produce non-digits as output in the absence of label information (see, e.g., Bortoli et al. (2021)). More generally, blending modes may or may not be desirable, since it can produce (e.g.) images very qualitatively different from those of the training distribution.

Other forms of generalization are possible. Factors we did not study, like learning dynamics, most likely also partly determine how diffusion models generalize. For example, the use of stochastic gradient descent introduces additional randomness that disfavors converging on sharp local optima (Smith & Le, 2018; Smith et al., 2020). It would be interesting to utilize recent theoretical tools (Bordelon & Pehlevan, 2023) to characterize how learning dynamics impacts generalization, especially in the rich (Geiger et al., 2020; Woodworth et al., 2020) rather than lazy learning regime.

Comment on memorization. Determining whether diffusion models memorize data (Somepalli et al., 2023a; Carlini et al., 2023), and if so how to address the issue (Vyas et al., 2023), has become a significant technical and societal issue. Our theory suggests that since generalization through variance happens primarily in boundary regions, diffusion models are unlikely to substantially generalize outliers. Since conditional models involve distributions of much higher effective dimension, one may expect that more training examples are ‘outlier-like’, and hence memorization should happen more often; this is consistent with the observations of Somepalli et al. (2023b). Our theory also suggests why duplications increase memorization: the existence of a strong boundary between modes, which requires modes to have comparable probability mass, is degraded.

Limitations of theoretical approach. Our theory is simplified in at least two ways. First, only a simple formulation of training (via DSM) and sampling (via the PF-ODE) from diffusion models is considered. There exist alternatives to DSM, like sliced score matching (Song et al., 2020), and alternative ways of sampling, including using auxiliary momentum-like variables (Dockhorn et al., 2022b). Also, our theoretical analysis neglects variation due to numerical integration schemes, even though these may matter in practice (Liu et al., 2022; Karras et al., 2022; Dockhorn et al., 2022a).

Second, we study only unconditional models for simplicity. This means that in particular do not consider diffusion coupled to attention layers, which enables the text-conditioning behind many of the most striking diffusion-model-related successes (Rombach et al., 2022; Blattmann et al., 2023).

Finally, we do not consider realistic architectures (like U-nets) and rich learning dynamics due to theoretical tractability. However, these challenges are not unique to the current setting. Despite our contribution’s simplicity, we hope that it nonetheless provides a foundation for others to more rigorously understand the inductive biases and generalization capabilities of diffusion models.

REFERENCES

- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/dbc4d84bfcfe2284ballbeffb853a8c4-Paper.pdf.
- Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/c4ef9c39b300931b69a36fb3dbb8d60e-Paper.pdf.
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 22563–22575, June 2023.
- Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution in wide neural networks. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 32240–32256. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/d027a5c93d484a4312cc486d399c62c1-Paper-Conference.pdf.
- Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution in wide neural networks*. *Journal of Statistical Mechanics: Theory and Experiment*, 2023(11): 114009, nov 2023. doi: 10.1088/1742-5468/ad01b0. URL <https://dx.doi.org/10.1088/1742-5468/ad01b0>.
- Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. Spectrum dependent learning curves in kernel regression and wide neural networks. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1024–1034. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/bordelon20a.html>.
- Valentin De Bortoli. Convergence of denoising diffusion models under the manifold hypothesis. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=MhK5aXo3gB>. Expert Certification.
- Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion schrödinger bridge with applications to score-based generative modeling. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=9BnCwiXB0ty>.
- Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks. *Nature Communications*, 12(1):2914, May 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-23103-1. URL <https://doi.org/10.1038/s41467-021-23103-1>.
- Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pp. 5253–5270, Anaheim, CA, August 2023. USENIX Association. ISBN 978-1-939133-37-3. URL <https://www.usenix.org/conference/usenixsecurity23/presentation/carlini>.
- Chen-Hao Chao, Wei-Fang Sun, Bo-Wun Cheng, Yi-Chen Lo, Chia-Che Chang, Yu-Lun Liu, Yu-Lin Chang, Chia-Ping Chen, and Chun-Yi Lee. Denoising likelihood score matching for conditional score-based data generation. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=LcF-EEt8cCC>.

- Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 4672–4712. PMLR, 23–29 Jul 2023a. URL <https://proceedings.mlr.press/v202/chen23o.html>.
- Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/69386f6bb1dfed68692a24c8686939b9-Paper.pdf.
- Sitan Chen, Sinho Chewi, Holden Lee, Yuanzhi Li, Jianfeng Lu, and Adil Salim. The probability flow ode is provably fast. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 6852–68575. Curran Associates, Inc., 2023b. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/d84a27ff694345aacc21c72097a69ea2-Paper-Conference.pdf.
- Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. In *The Eleventh International Conference on Learning Representations*, 2023c. URL https://openreview.net/forum?id=zyLVMgsZ0U_.
- Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/ae614c557843b1df326cb29c57225459-Paper.pdf.
- Taco Cohen and Max Welling. Group equivariant convolutional networks. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 2990–2999, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/cohen16.html>.
- A. Crisanti and H. Sompolinsky. Path integral approach to random neural networks. *Phys. Rev. E*, 98:062120, Dec 2018. doi: 10.1103/PhysRevE.98.062120. URL <https://link.aps.org/doi/10.1103/PhysRevE.98.062120>.
- Salman Ul Hassan Dar, Arman Ghanaat, Jannik Kahmann, Isabelle Ayx, Theano Papavassiliu, Stefan O. Schoenberg, and Sandy Engelhardt. Investigating data memorization in 3d latent diffusion models for medical image synthesis. In *Deep Generative Models: Third MICCAI Workshop, DGM4MICCAI 2023, Held in Conjunction with MICCAI 2023, Vancouver, BC, Canada, October 8, 2023, Proceedings*, pp. 56–65, Berlin, Heidelberg, 2024. Springer-Verlag. ISBN 978-3-031-53766-0. doi: 10.1007/978-3-031-53767-7_6. URL https://doi.org/10.1007/978-3-031-53767-7_6.
- Tim Dockhorn, Arash Vahdat, and Karsten Kreis. GENIE: Higher-order denoising diffusion solvers. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022a. URL <https://openreview.net/forum?id=LKEYuYNOqx>.
- Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Score-based generative modeling with critically-damped langevin diffusion. In *International Conference on Learning Representations*, 2022b. URL <https://openreview.net/forum?id=CzceR82CYc>.
- Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2020 (11):113301, nov 2020. doi: 10.1088/1742-5468/abc4de. URL <https://dx.doi.org/10.1088/1742-5468/abc4de>.

- Yinbin Han, Meisam Razaviyayn, and Renyuan Xu. Neural network-based score estimation in diffusion models: Optimization and generalization. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=h8GeqOxt4d>.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf.
- Zahra Kadkhodaie, Florentin Guth, Eero P Simoncelli, and Stéphane Mallat. Generalization in diffusion models arises from geometry-adaptive harmonic representations. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ANvmVS2Yr0>.
- Mason Kamb and Surya Ganguli. An analytic theory of creativity in convolutional diffusion models. *arXiv e-prints*, art. arXiv:2412.20292, December 2024. doi: 10.48550/arXiv.2412.20292.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=k7FuTOWMOc7>.
- Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing and improving the training dynamics of diffusion models, 2024.
- Hagen Kleinert. *Path integrals in quantum mechanics, statistics, polymer physics, and financial markets*. World Scientific Publishing Company, 2006.
- Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778*, 2022.
- Lorenzo Luzzi, Paul M Mayer, Josue Casco-Rodriguez, Ali Siahkoohi, and Richard Baraniuk. Boomerang: Local sampling on image manifolds using diffusion models. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=NYdThkjNW1>.
- P. C. Martin, E. D. Siggia, and H. A. Rose. Statistical dynamics of classical systems. *Phys. Rev. A*, 8:423–437, Jul 1973. doi: 10.1103/PhysRevA.8.423. URL <https://link.aps.org/doi/10.1103/PhysRevA.8.423>.
- Francesca Mignacco, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborová. Dynamical mean-field theory for stochastic gradient descent in gaussian mixture classification. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 9540–9550. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/6c81c83c4bd0b58850495f603ab45a93-Paper.pdf.
- Matthew Niedoba, Berend Zwartsenberg, Kevin Murphy, and Frank Wood. Towards a mechanistic explanation of diffusion model generalization. *arXiv e-prints*, art. arXiv:2411.19339, February 2025. URL <https://arxiv.org/abs/2411.19339>.
- Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5301–5310. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/rahaman19a.html>.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn 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.

- Haozhe Shan and Blake Bordelon. A theory of neural tangent kernel alignment and its influence on training, 2022. URL <https://arxiv.org/abs/2105.14301>.
- Samuel Smith, Erich Elsen, and Soham De. On the generalization benefit of noise in stochastic gradient descent. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 9058–9067. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/smith20a.html>.
- Samuel L. Smith and Quoc V. Le. A bayesian perspective on generalization and stochastic gradient descent. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=BJij4yg0Z>.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2256–2265, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6048–6058, June 2023a.
- Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Understanding and mitigating copying in diffusion models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b. URL <https://openreview.net/forum?id=HtMXRGbUMt>.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/3001ef257407d5a371a96dcd947c7d93-Paper.pdf>.
- Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. In Ryan P. Adams and Vibhav Gogate (eds.), *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, volume 115 of *Proceedings of Machine Learning Research*, pp. 574–584. PMLR, 22–25 Jul 2020. URL <https://proceedings.mlr.press/v115/song20a.html>.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=PXTIG12RRHS>.
- George Stein, Jesse C. Cresswell, Rasa Hosseinzadeh, Yi Sui, Brendan Leigh Ross, Valentin Vilecroze, Zhaoyan Liu, Anthony L. Caterini, Eric Taylor, and Gabriel Loaiza-Ganem. Exposing flaws of generative model evaluation metrics and their unfair treatment of diffusion models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=08zf7kTOoh>.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011. doi: 10.1162/NECO_a.00142.
- Nikhil Vyas, Sham M. Kakade, and Boaz Barak. On provable copyright protection for generative models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 35277–35299. PMLR, 2023. URL <https://proceedings.mlr.press/v202/vyas23b.html>.

- Binxu Wang and John Vastola. The unreasonable effectiveness of gaussian score approximation for diffusion models and its applications. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=I0uknSHM2j>.
- Binxu Wang and John J. Vastola. Diffusion models generate images like painters: an analytical theory of outline first, details later. *arXiv e-prints*, art. arxiv:2303.02490, March 2023. URL <https://arxiv.org/abs/2303.02490>.
- Blake Woodworth, Suriya Gunasekar, Jason D. Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In Jacob Abernethy and Shivani Agarwal (eds.), *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pp. 3635–3673. PMLR, 09–12 Jul 2020. URL <https://proceedings.mlr.press/v125/woodworth20a.html>.
- Yilun Xu, Shangyuan Tong, and Tommi S. Jaakkola. Stable target field for reduced variance score estimation in diffusion models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=WmIwYTd0YTF>.
- Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Comput. Surv.*, 56(4), nov 2023. ISSN 0360-0300. doi: 10.1145/3626235. URL <https://doi.org/10.1145/3626235>.
- TaeHo Yoon, Joo Young Choi, Sehyun Kwon, and Ernest K. Ryu. Diffusion probabilistic models generalize when they fail to memorize. In *ICML 2023 Workshop on Structured Probabilistic Inference & Generative Modeling*, 2023. URL <https://openreview.net/forum?id=shciCbSk9h>.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Sy8gdB9xx>.
- Huijie Zhang, Jinfan Zhou, Yifu Lu, Minzhe Guo, Peng Wang, Liyue Shen, and Qing Qu. The emergence of reproducibility and consistency in diffusion models. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=HsliOqZkc0>.

Appendix

See <https://github.com/john-vastola/gtv-iclr25> for code that produces Fig. 1-3.

Table of Contents

A Optimizing objective reproduces training distribution	16
A.1 Denoising score matching preserves optima of naive objective	16
A.2 Training distribution reproduction	17
B Covariance of proxy score	18
B.1 Computing covariance of proxy score	18
B.2 Connection to Fisher information	18
B.3 Explicit covariance for isotropic Gaussian mixture training distribution	18
C Boundary regions: definition and Bayesian interpretation	20
D Path-integral representation of learned distribution	21
D.1 Warm-up: Deriving a path-integral representation of the PF-ODE	21
D.2 Deriving a path-integral representation of a more general process	22
D.3 Averaging learned distribution over sample realizations	23
E Naive score estimators generalize: details	24
F Linear score estimator: details	25
F.1 Definition of linear score model	25
F.2 Optimum of DSM objective for linear score model	25
F.3 Optimum of DSM objective given a finite number of samples	26
F.4 Linear score model estimator is unbiased	27
F.5 Computing the V-kernel of the linear score model	28
G Neural network score estimator in NTK regime: details	29
G.1 Definition of neural network model	29
G.2 Learned score after full-batch gradient descent	29
G.3 Computing the V-kernel of the NTK model	31
H Results for noise prediction formulation	33
I Results for denoiser formulation	34
J Benign properties of generalization through variance	35
J.1 Single points are not generalized	35
J.2 Dimensionality of data distribution is preserved	35
J.3 Variance is not added far from data distribution examples	35
J.4 Following average score field is most likely	36
J.5 Training data are more likely to be sampled when noise is small	36
K Memorization and the V-kernel in the small noise limit	37
K.1 Setting up the semiclassical approximation	37
K.2 Semiclassical approximation of the learned distribution	38
K.3 Quantifying memorization in the semiclassical regime	39

A OPTIMIZING OBJECTIVE REPRODUCES TRAINING DISTRIBUTION

In this appendix, we characterize the optima of the naive and DSM objectives introduced in Sec. 2, and in particular show that one (naively) theoretically expects diffusion models to reproduce the training distribution in the absence of expressivity-related constraints.

A.1 DENOISING SCORE MATCHING PRESERVES OPTIMA OF NAIVE OBJECTIVE

First, we reestablish the well-known fact that the optima of the naive objective

$$J_0(\boldsymbol{\theta}) := \frac{1}{2} \mathbb{E}_{t,\mathbf{x}} \{ \lambda_t \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{s}(\mathbf{x}, t)\|_2^2 \} = \int \frac{\lambda_t}{2} \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{s}(\mathbf{x}, t)\|_2^2 p(\mathbf{x}|t)p(t) d\mathbf{x}dt \quad (19)$$

and DSM objective

$$\begin{aligned} J_1(\boldsymbol{\theta}) &:= \frac{1}{2} \mathbb{E}_{t,\mathbf{x}_0,\mathbf{x}} \{ \lambda_t \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)\|_2^2 \} \\ &= \int \frac{\lambda_t}{2} \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)\|_2^2 p(\mathbf{x}|\mathbf{x}_0, t)p_{data}(\mathbf{x}_0)p(t) d\mathbf{x}d\mathbf{x}_0dt \end{aligned} \quad (20)$$

are the same (Vincent, 2011; Song & Ermon, 2019; Song et al., 2021). Assume that $\mathbf{x}, \mathbf{x}_0 \in \mathbb{R}^D$ and that $\boldsymbol{\theta} \in \mathbb{R}^F$. The gradient of J_0 with respect to $\boldsymbol{\theta}$ is

$$\frac{\partial J_0}{\partial \boldsymbol{\theta}} = \int \lambda_t \frac{\partial \hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t)^T}{\partial \boldsymbol{\theta}} [\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{s}(\mathbf{x}, t)] p(\mathbf{x}|t)p(t) d\mathbf{x}dt \quad (21)$$

where $\partial \hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t)/\partial \boldsymbol{\theta}$ is the $D \times F$ Jacobian matrix of the score estimator. The gradient of J_1 is

$$\frac{\partial J_1}{\partial \boldsymbol{\theta}} = \int \lambda_t \frac{\partial \hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t)^T}{\partial \boldsymbol{\theta}} [\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)] p(\mathbf{x}|\mathbf{x}_0, t)p_{data}(\mathbf{x}_0)p(t) d\mathbf{x}d\mathbf{x}_0dt. \quad (22)$$

At this point, we make two observations about the gradient of J_1 . First, the term on the left does not depend on \mathbf{x}_0 , so we can marginalize over \mathbf{x}_0 . Explicitly,

$$\int \lambda_t \frac{\partial \hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t)^T}{\partial \boldsymbol{\theta}} \hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) p(\mathbf{x}|\mathbf{x}_0, t)p_{data}(\mathbf{x}_0)p(t) d\mathbf{x}d\mathbf{x}_0dt = \int \lambda_t \frac{\partial \hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t)^T}{\partial \boldsymbol{\theta}} \hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) p(\mathbf{x}|t)p(t) d\mathbf{x}dt.$$

Second, the term on the right only depends on \mathbf{x}_0 through the proxy score target. Moreover,

$$\begin{aligned} \int \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0) p(\mathbf{x}|\mathbf{x}_0, t)p_{data}(\mathbf{x}_0) d\mathbf{x}_0 &= \int \nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{x}_0, t) p(\mathbf{x}|\mathbf{x}_0, t)p_{data}(\mathbf{x}_0) d\mathbf{x}_0 \\ &= \int \nabla_{\mathbf{x}} p(\mathbf{x}|\mathbf{x}_0, t)p_{data}(\mathbf{x}_0) d\mathbf{x}_0 \\ &= \nabla_{\mathbf{x}} \int p(\mathbf{x}|\mathbf{x}_0, t)p_{data}(\mathbf{x}_0) d\mathbf{x}_0 \\ &= \nabla_{\mathbf{x}} p(\mathbf{x}|t) \\ &= \mathbf{s}(\mathbf{x}, t)p(\mathbf{x}|t). \end{aligned} \quad (23)$$

Hence, the gradient of J_0 is the same as the gradient of J_1 , so they have the same optima. If the score approximator is arbitrarily expressive and smooth in its parameters, we in particular have that the true score (a global minimum of J_0) is an optimum of the DSM objective.

This optimum is *also* the global minimum of J_1 . Note that J_1 can be written as

$$\begin{aligned} &\mathbb{E}_{t,\mathbf{x}_0,\mathbf{x}} \left\{ \frac{\lambda_t}{2} \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{s}(\mathbf{x}, t) + \mathbf{s}(\mathbf{x}, t) - \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)\|_2^2 \right\} \\ &= \mathbb{E}_{t,\mathbf{x}_0,\mathbf{x}} \left\{ \frac{\lambda_t}{2} \left(\|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{s}(\mathbf{x}, t)\|_2^2 + 2[\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{s}(\mathbf{x}, t)] \cdot [\mathbf{s}(\mathbf{x}, t) - \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)] + \|\mathbf{s}(\mathbf{x}, t) - \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)\|_2^2 \right) \right\}. \end{aligned}$$

The first term is precisely equal to J_0 . The second term vanishes, since (as shown by Eq. 23)

$$\mathbb{E}_{\mathbf{x}_0|\mathbf{x},t}[\tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)] = \mathbf{s}(\mathbf{x}, t). \quad (24)$$

Hence, we have that

$$J_1 = J_0 + \frac{1}{2} \mathbb{E}_{t,\mathbf{x}} \{ \lambda_t \text{tr}(\text{Cov}_{\mathbf{x}_0|\mathbf{x},t}(\tilde{\mathbf{s}})) \}. \quad (25)$$

In words: J_1 is equal to J_0 up to a $\boldsymbol{\theta}$ -independent term that is a weighted combination of proxy score variances.

A.2 TRAINING DISTRIBUTION REPRODUCTION

In practice, the training set consists of $1 \leq M < \infty$ examples (e.g., images) which together define

$$p_{data}(\mathbf{x}_0) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x}_0 - \boldsymbol{\mu}_m). \quad (26)$$

The corresponding ‘corrupted’ distribution, given our choice of forward process (see Sec. 2), is

$$p(\mathbf{x}|t) = \frac{1}{M} \sum_{m=1}^M \mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_m, \mathbf{S}_t). \quad (27)$$

Usually, model updates utilize batches of samples from $p(\mathbf{x}, \mathbf{x}_0, t)$ (Song et al., 2021; Karras et al., 2022). As training proceeds, the model sees an ever larger number P of samples from this distribution, making the empirical objective

$$J_1(\boldsymbol{\theta}; P) := \frac{1}{P} \sum_{n=1}^P \frac{\lambda(t^{(n)})}{2} \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}^{(n)}, t^{(n)}) - \tilde{\mathbf{s}}(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}_0^{(n)})\|_2^2, \quad (28)$$

where the n superscripts index different (independent) samples from $p(\mathbf{x}, \mathbf{x}_0, t) = p(\mathbf{x}|\mathbf{x}_0, t)p_{data}(\mathbf{x}_0)p(t)$. For P sufficiently large, by the central limit theorem, we expect the empirical objective to be extremely close to the true objective, and hence share its global minimum. But the global minimum is the true score, i.e.,

$$\mathbf{s}(\mathbf{x}, t) = \sum_{m=1}^M \mathbf{S}_t^{-1}(\alpha_t \boldsymbol{\mu}_m - \mathbf{x}) \frac{\mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_m, \mathbf{S}_t)}{\sum_{m'} \mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_{m'}, \mathbf{S}_t)}.$$

Since integrating the PF-ODE using this score produces samples from $p_{data}(\mathbf{x}_0)$ —as $t \rightarrow 0$, $\mathbf{S}_t \rightarrow \mathbf{0}_D$, so the asymptotic ‘force’ pushing \mathbf{x}_t towards an example becomes infinitely strong—we expect expressive diffusion models trained on the DSM objective using a large number of samples to reproduce training examples.

B COVARIANCE OF PROXY SCORE

In this appendix, we compute the covariance of the proxy score $\tilde{s}(\mathbf{x}, t; \mathbf{x}_0) := \nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{x}_0, t)$ with respect to $p(\mathbf{x}_0|\mathbf{x}, t)$. We also show how this covariance is connected to Fisher information, and explicitly compute it in the case that $p_{data}(\mathbf{x}_0)$ is an isotropic Gaussian mixture.

B.1 COMPUTING COVARIANCE OF PROXY SCORE

Note that

$$\frac{\partial^2}{\partial x_i \partial x_j} p(\mathbf{x}|\mathbf{x}_0, t) = [-S_{t,ij}^{-1} + \tilde{s}_i \tilde{s}_j] p(\mathbf{x}|\mathbf{x}_0, t). \quad (29)$$

Using this fact, we can write

$$\begin{aligned} \text{Cov}_{\mathbf{x}_0|\mathbf{x}, t}(\tilde{s}_i, \tilde{s}_j) &= \int \tilde{s}_i \tilde{s}_j \frac{p(\mathbf{x}|\mathbf{x}_0, t) p_{data}(\mathbf{x}_0)}{p(\mathbf{x}|t)} d\mathbf{x}_0 - s_i s_j \\ &= \int \frac{1}{p(\mathbf{x}|t)} \left[S_{t,ij}^{-1} + \frac{\partial^2}{\partial x_i \partial x_j} \right] p(\mathbf{x}|\mathbf{x}_0, t) p_{data}(\mathbf{x}_0) d\mathbf{x}_0 - s_i s_j \\ &= \int \frac{1}{p(\mathbf{x}|t)} \left[S_{t,ij}^{-1} + \frac{\partial^2}{\partial x_i \partial x_j} \right] p(\mathbf{x}|t) \frac{p(\mathbf{x}|\mathbf{x}_0, t) p_{data}(\mathbf{x}_0)}{p(\mathbf{x}|t)} d\mathbf{x}_0 - s_i s_j \quad (30) \\ &= S_{t,ij}^{-1} + \frac{1}{p(\mathbf{x}|t)} \frac{\partial^2 p(\mathbf{x}|t)}{\partial x_i \partial x_j} - s_i s_j \\ &= S_{t,ij}^{-1} + \frac{\partial^2}{\partial x_i \partial x_j} \log p(\mathbf{x}|t). \end{aligned}$$

B.2 CONNECTION TO FISHER INFORMATION

By definition, if $p(\mathbf{x}_0|\mathbf{x}, t)$ is viewed as a distribution with parameter vector \mathbf{x} , and t is viewed as a hyperparameter, the Fisher information \mathcal{I}_F is defined as

$$\begin{aligned} \mathcal{I}_F(\mathbf{x}|t) &:= \int \frac{\partial \log p(\mathbf{x}_0|\mathbf{x}, t)}{\partial x_i} \cdot \frac{\partial \log p(\mathbf{x}_0|\mathbf{x}, t)}{\partial x_j} p(\mathbf{x}_0|\mathbf{x}, t) d\mathbf{x}_0 \\ &= \int \left[\frac{\partial \log p(\mathbf{x}|\mathbf{x}_0, t)}{\partial x_i} - \frac{\partial \log p(\mathbf{x}|t)}{\partial x_i} \right] \left[\frac{\partial \log p(\mathbf{x}|\mathbf{x}_0, t)}{\partial x_j} - \frac{\partial \log p(\mathbf{x}|t)}{\partial x_j} \right] p(\mathbf{x}_0|\mathbf{x}, t) d\mathbf{x}_0 \quad (31) \\ &= \int [\tilde{s}_i - s_i] [\tilde{s}_j - s_j] p(\mathbf{x}_0|\mathbf{x}, t) d\mathbf{x}_0 \\ &= \text{Cov}_{\mathbf{x}_0|\mathbf{x}, t}(\tilde{s}_i, \tilde{s}_j). \end{aligned}$$

B.3 EXPLICIT COVARIANCE FOR ISOTROPIC GAUSSIAN MIXTURE TRAINING DISTRIBUTION

Suppose that $p(\mathbf{x}_0)$ and $p(\mathbf{x}|t)$ are

$$p_{data}(\mathbf{x}_0) = \frac{1}{M} \sum_m \mathcal{N}(\mathbf{x}_0; \boldsymbol{\mu}_m, \sigma_0^2 \mathbf{I}) \quad p(\mathbf{x}|t) = \frac{1}{M} \sum_m \mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_m, \alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t). \quad (32)$$

Note that the delta mixture case is an example ($\sigma_0^2 = 0$). Define the softmax distribution

$$p(m|\mathbf{x}, t) := \frac{\mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_m, \alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)}{\sum_{m'} \mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_{m'}, \alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)} \quad (33)$$

on $\mathcal{M} = \{1, \dots, M\}$. This distribution, whose moments determine the proxy score covariance, has a Bayesian interpretation: it corresponds to an ideal observer's belief about the outcome x_0 , given that said observer is in state x at time t .

The first and second derivatives of $p(\mathbf{x}|t)$ can be written in terms of expectations with respect to this distribution, since

$$\frac{1}{p(\mathbf{x}|t)} \frac{\partial p(\mathbf{x}|t)}{\partial \mathbf{x}} = \sum_m (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} (\alpha_t \boldsymbol{\mu}_m - \mathbf{x}) p(m|\mathbf{x}, t) = (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} (\alpha_t \langle \boldsymbol{\mu} \rangle_{\mathcal{M}} - \mathbf{x})$$

and the Hessian matrix ($H_{ij} := \partial_{ij}^2 p(\mathbf{x}|t)$) is

$$\begin{aligned} \frac{\mathbf{H}}{p(\mathbf{x}|t)} &= \sum_m [-(\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} + (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} (\alpha_t \boldsymbol{\mu}_m - \mathbf{x})(\alpha_t \boldsymbol{\mu}_m - \mathbf{x})^T (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1}] p(m|\mathbf{x}, t) \\ &= -(\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} + (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} \mathbb{E}_{\mathcal{M}} \{ (\alpha_t \boldsymbol{\mu} - \mathbf{x})(\alpha_t \boldsymbol{\mu} - \mathbf{x})^T \} (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} \\ &= -(\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} + (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} [\alpha_t^2 \text{Cov}_{\mathcal{M}}(\boldsymbol{\mu}) + (\alpha_t \langle \boldsymbol{\mu} \rangle_{\mathcal{M}} - \mathbf{x})(\alpha_t \langle \boldsymbol{\mu} \rangle_{\mathcal{M}} - \mathbf{x})^T] (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1}. \end{aligned}$$

Then we have

$$\frac{\partial^2 \log p(\mathbf{x}|t)}{\partial x_i \partial x_j} = -(\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} + \alpha_t^2 (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} \text{Cov}_{\mathcal{M}}(\boldsymbol{\mu}) (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} \quad (34)$$

and hence that

$$\text{Cov}_{\mathbf{x}_0|\mathbf{x},t}(\tilde{\mathbf{s}}) = \mathbf{S}_t^{-1} - (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} + \alpha_t^2 (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1} \text{Cov}_{\mathcal{M}}(\boldsymbol{\mu}) (\alpha_t^2 \sigma_0^2 \mathbf{I} + \mathbf{S}_t)^{-1}.$$

For a delta mixture training distribution, since $\sigma_0^2 = 0$, the covariance simplifies to

$$\text{Cov}_{\mathbf{x}_0|\mathbf{x},t}(\tilde{\mathbf{s}}) = \alpha_t^2 \mathbf{S}_t^{-1} \text{Cov}_{\mathcal{M}}(\boldsymbol{\mu}) \mathbf{S}_t^{-1}. \quad (35)$$

The above equation implies that the covariance of the proxy score is, up to scaling, the same as uncertainty about \mathbf{x}_0 given \mathbf{x} and t .

C BOUNDARY REGIONS: DEFINITION AND BAYESIAN INTERPRETATION

A key concept used throughout this paper is that of a *boundary region*, which we informally define as a region of \mathbb{R}^D between two or more training examples in the case that p_{data} is discrete. Intuitively, these regions correspond to the ‘gaps’ in the training distribution, and a reasonable generalization strategy is to fill them in.

In this appendix, we briefly comment that the notion of a boundary region can be made more precise via Bayes’ theorem: for a given reverse diffusion time t , boundary regions are sets of x values for which uncertainty about the endpoint x_0 is particularly high. For discrete data distributions, such regions coincide with sets of states between training examples, since one is maximally uncertain about the endpoint when x is equidistant from two or more training examples.

An illustrative one-dimensional example involves two training examples at $x_0 = \pm\mu$. The noise-corrupted data distribution at time t is

$$p(x|t) = \frac{1}{2}\mathcal{N}(x; \alpha_t\mu, \sigma_t^2) + \frac{1}{2}\mathcal{N}(x; -\alpha_t\mu, \sigma_t^2), \quad (36)$$

so the posterior endpoint estimate given x, t is the softmax distribution (see also Appendix B)

$$\begin{aligned} p(x_0|x, t) &= \frac{p(x|x_0, t)p_{data}(x_0)}{\sum_{x_0} p(x|x_0, t)p_{data}(x_0)} \\ &= \frac{\mathcal{N}(x; \alpha_t\mu, \sigma_t^2)}{\mathcal{N}(x; \alpha_t\mu, \sigma_t^2) + \mathcal{N}(x; -\alpha_t\mu, \sigma_t^2)}\delta(x_0 - \mu) + \frac{\mathcal{N}(x; -\alpha_t\mu, \sigma_t^2)}{\mathcal{N}(x; \alpha_t\mu, \sigma_t^2) + \mathcal{N}(x; -\alpha_t\mu, \sigma_t^2)}\delta(x_0 + \mu). \end{aligned}$$

The mean $\mathbb{E}[x_0|x, t]$ of this distribution is

$$\mathbb{E}[x_0|x, t] = \sum_{x_0} x_0 p(x_0|x, t) = \mu \tanh\left(\frac{\alpha_t\mu}{\sigma_t^2}x\right). \quad (37)$$

The interpretation of this quantity is interesting in light of the ‘Bayesian guessing game’ metaphor for score learning (see, e.g., Kamb & Ganguli (2024)). One imagines that one starts at an unknown x_0 (here, either $+\mu$ or $-\mu$), and then noise is added according to the forward process until time t . Given that an observer is in state x at time t , what was the likely starting point x_0 ? The quantity $\mathbb{E}[x_0|x, t]$ is the Bayes-optimal solution to this problem.

In the context of this toy example, it has the following form. If x is very positive, one tends to believe $x_0 = +\mu$; if x is very negative, one tends to believe $x_0 = -\mu$. For intermediate x , especially near $x = 0$, uncertainty is highest, and $\mathbb{E}[x_0|x, t]$ is near zero, since the observer could have plausibly started at either $x_0 = +\mu$ or $x_0 = -\mu$.

This high uncertainty allows us to formalize the idea that the region between $+\mu$ and $-\mu$, especially near $x = 0$, is a boundary region. Quantitatively, we have

$$\text{var}(x_0|x, t) = \sum_{x_0} x_0^2 p(x_0|x, t) - \mathbb{E}[x_0|x, t]^2 = \mu^2 \left[1 - \tanh^2\left(\frac{\alpha_t\mu}{\sigma_t^2}x\right) \right] = \frac{\mu^2}{\cosh^2\left(\frac{\alpha_t\mu}{\sigma_t^2}x\right)} \quad (38)$$

or equivalently

$$\sqrt{\text{var}(x_0|x, t)} = \frac{\mu}{\cosh\left(\frac{\alpha_t\mu}{\sigma_t^2}x\right)}. \quad (39)$$

At $x = 0$, the standard deviation of $p(x_0|x, t)$ equals μ , i.e., there is maximum uncertainty about the starting point x_0 . Moreover, it is fairly high until $x \approx \frac{\sigma_t^2}{\alpha_t\mu}$, which also shows that the effective size of a boundary region is smaller at smaller noise scales. Said differently, the basins of attraction surrounding each training example become increasingly sharp as $\sigma_t \rightarrow 0$.

D PATH-INTEGRAL REPRESENTATION OF LEARNED DISTRIBUTION

In this appendix, we derive a path-integral description of the ‘typical’ distribution learned by diffusion models. We do this in three stages. First, we derive a path-integral description of the PF-ODE. Next, we derive a path-integral description of a more general kind of stochastic process. Finally, we show that averaging the path-integral representation of the PF-ODE over sample realizations produces a path integral whose dynamics correspond to those of the aforementioned stochastic process.

D.1 WARM-UP: DERIVING A PATH-INTEGRAL REPRESENTATION OF THE PF-ODE

A general ODE can be written as

$$\dot{\mathbf{x}}_t = \mathbf{f}(\mathbf{x}_t, t) \quad (40)$$

where $\mathbf{x}_t \in \mathbb{R}^D$ and $t \in [\epsilon, T]$. We will assume that \mathbf{f} is smooth to avoid technical issues. If we discretize time, and slightly abuse notation by using t and T to refer to integer-valued indices instead of real-valued times, we can write the trajectory as $\{\mathbf{x}_T, \mathbf{x}_{T-1}, \dots, \mathbf{x}_1, \mathbf{x}_0\}$ and the corresponding updates in the form

$$\mathbf{x}_t = \mathbf{x}_{t+1} - \mathbf{f}(\mathbf{x}_{t+1}, t+1)\Delta t. \quad (41)$$

Note that our discretization corresponds to a first-order Euler update scheme. In the small Δt limit, this specific choice does not matter, even if it matters in practice; we use it to slightly simplify our argument. Conditional on the initial point \mathbf{x}_T , the probability of reaching another point \mathbf{x}_0 after T backwards-time steps is

$$p(\mathbf{x}_0|\mathbf{x}_T) = \int \delta(\mathbf{x}_0 - \mathbf{x}_1 + \mathbf{f}(\mathbf{x}_1, 1)\Delta t) \cdots \delta(\mathbf{x}_{T-1} - \mathbf{x}_T + \mathbf{f}(\mathbf{x}_T, T)\Delta t) d\mathbf{x}_1 \cdots d\mathbf{x}_{T-1} \quad (42)$$

where δ is the Dirac delta function. Here, we will employ a well-known integral representation of the Dirac delta function:

$$\delta(\mathbf{x} - \mathbf{x}') = \int \frac{d\mathbf{p}}{(2\pi)^D} \exp\{-i\mathbf{p} \cdot (\mathbf{x} - \mathbf{x}')\} \quad (43)$$

where \mathbf{p} is integrated over all of \mathbb{R}^D . Our expression for $p(\mathbf{x}_0|\mathbf{x}_T)$ becomes

$$p(\mathbf{x}_0|\mathbf{x}_T) = \int \frac{d\mathbf{p}_0}{(2\pi)^D} \frac{d\mathbf{x}_1 d\mathbf{p}_1}{(2\pi)^D} \cdots \frac{d\mathbf{x}_{T-1} d\mathbf{p}_{T-1}}{(2\pi)^D} \exp\left\{\sum_{t=0}^{T-1} -i\mathbf{p}_t \cdot [\mathbf{x}_t - \mathbf{x}_{t+1} + \mathbf{f}(\mathbf{x}_{t+1}, t+1)\Delta t]\right\}. \quad (44)$$

Schematically, we can write this path integral as a ‘sum over paths’

$$p(\mathbf{x}_0|\mathbf{x}_T) = \int \mathcal{D}[\mathbf{p}_t] \mathcal{D}[\mathbf{x}_t] \exp\left\{\int_{\epsilon}^T -i\mathbf{p}_t \cdot [-\dot{\mathbf{x}}_t + \mathbf{f}(\mathbf{x}_t, t)] dt\right\}, \quad (45)$$

although explicitly using this form is unnecessary for our purposes. (This is good, since remaining in discrete time allows us to avoid various thorny mathematical issues.) For the particular choice of \mathbf{f} associated with the PF-ODE, we have discrete and schematic forms

$$p(\mathbf{x}_0|\mathbf{x}_T) = \int \frac{d\mathbf{p}_0}{(2\pi)^D} \frac{d\mathbf{x}_1 d\mathbf{p}_1}{(2\pi)^D} \cdots \frac{d\mathbf{x}_{T-1} d\mathbf{p}_{T-1}}{(2\pi)^D} e^{\sum_{t=0}^{T-1} -i\mathbf{p}_t \cdot [\mathbf{x}_t - \mathbf{x}_{t+1} - (\beta_{t+1}\mathbf{x}_{t+1} + \mathbf{D}_{t+1}\mathbf{s}(\mathbf{x}_{t+1}, t+1))\Delta t]}$$

$$p(\mathbf{x}_0|\mathbf{x}_T) = \int \mathcal{D}[\mathbf{p}_t] \mathcal{D}[\mathbf{x}_t] \exp\left\{\int_{\epsilon}^T i\mathbf{p}_t \cdot [\dot{\mathbf{x}}_t + \beta_t \mathbf{x}_t + \mathbf{D}_t \mathbf{s}(\mathbf{x}_t, t)] dt\right\}.$$

D.2 DERIVING A PATH-INTEGRAL REPRESENTATION OF A MORE GENERAL PROCESS

Consider a more general type of backwards, discrete-time stochastic process. Once again, suppose that a variable $\mathbf{x}_t \in \mathbb{R}^D$ evolves backwards in time from an initial point \mathbf{x}_T . But this time, suppose that the transition between \mathbf{x}_{t+1} and \mathbf{x}_t depends upon some set of K independent standard normal random variables $\{\xi_k\}$. In particular, suppose that discrete-time updates have the form

$$x_{t,j} = x_{t+1,j} - f_j(\mathbf{x}_{t+1}, t+1)\Delta t + \sum_{k=1}^K G_{jk}(\mathbf{x}_{t+1}, t+1) \xi_k \Delta t, \quad (46)$$

i.e., updates are the same as before except for the new noise term. In general, the noise term is quite complicated; \mathbf{G} is a $D \times K$ matrix which can depend explicitly on both the current state and the current time. The process described by the above updates is generally not Markov, since noise added at different time steps can depend on some of the same ξ_k variables, and hence the amount of noise added at one time step can be correlated with the amount of noise added at some other time step.

What is the distribution of \mathbf{x}_0 , the result of T steps of this process, conditional on a starting point \mathbf{x}_T ? We know that each update depends only on the previous state and the noise variables, so

$$p(\mathbf{x}_0|\mathbf{x}_T) = \int p(\mathbf{x}_0|\mathbf{x}_1, \{\xi_k\})p(\mathbf{x}_1|\mathbf{x}_2, \{\xi_k\}) \cdots p(\mathbf{x}_{T-1}|\mathbf{x}_T, \{\xi_k\}) p(\{\xi_k\}) d\mathbf{x}_1 \cdots d\mathbf{x}_{T-1} d\{\xi_k\}.$$

In particular, conditional on the previous state and the noise variables, updates are deterministic. This allows us to write the above transition probability as

$$\int \left[\prod_{j=1}^D \prod_{t=0}^{T-1} \delta \left(x_{t,j} - x_{t+1,j} + f_j(\mathbf{x}_{t+1}, t+1)\Delta t + \sum_{k=1}^K G_{jk}(\mathbf{x}_{t+1}, t+1) \xi_k \Delta t \right) \right] p(\{\xi_k\}) d\mathbf{x}_1 \cdots d\mathbf{x}_{T-1} d\{\xi_k\}.$$

Using the same integral representation of the Dirac delta function that we used above, this becomes

$$\int e^{\sum_{t,j} -ip_{t,j} [x_{t,j} - x_{t+1,j} + f_j(\mathbf{x}_{t+1}, t+1)\Delta t + \sum_{k=1}^K G_{jk}(\mathbf{x}_{t+1}, t+1) \xi_k \Delta t]} p(\{\xi_k\}) \frac{d\mathbf{p}_0}{(2\pi)^D} \frac{d\mathbf{x}_1 d\mathbf{p}_1}{(2\pi)^D} \cdots \frac{d\mathbf{x}_{T-1} d\mathbf{p}_{T-1}}{(2\pi)^D} d\{\xi_k\}.$$

Although this appears to be extremely complicated, it can be considerably simplified by doing the integral over the noise variables. Since the noise variables are all independent and standard normal,

$$p(\{\xi_k\}) = \frac{1}{(2\pi)^{k/2}} \exp \left\{ -\frac{\xi_1^2}{2} - \cdots - \frac{\xi_K^2}{2} \right\}. \quad (47)$$

Hence, the integral over the noise variables is a typical Gaussian integral with a linear term. We can save time by recognizing the integral as essentially computing the characteristic function of a standard normal; more precisely, we have

$$\begin{aligned} I_k &= \int \exp \left\{ -i\xi_k \sum_{t=0}^{T-1} \sum_{j=1}^D p_{t,j} G_{jk}(\mathbf{x}_{t+1}, t+1)\Delta t \right\} \frac{e^{-\xi_k^2/2}}{\sqrt{2\pi}} d\xi_k \\ &= \exp \left\{ -\frac{1}{2} \sum_{t=0}^{T-1} \sum_{t'=0}^{T-1} \sum_{j=1}^D \sum_{j'=1}^D p_{t,j} G_{jk}(\mathbf{x}_{t+1}, t+1) G_{j'k}(\mathbf{x}_{t'+1}, t'+1) p_{t',j'} \Delta t \Delta t \right\} \end{aligned} \quad (48)$$

for each ξ_k . Putting everything together, we find that $p(\mathbf{x}_0|\mathbf{x}_T)$ can be written

$$\int e^{\sum_{t,j} -ip_{t,j} [x_{t,j} - x_{t+1,j} + f_j(\mathbf{x}_{t+1}, t+1)\Delta t] - \frac{1}{2} \sum_{t,t',j,j'} \sum_{k=1}^K p_{t,j} G_{jk}(\mathbf{x}_{t+1}, t+1) G_{j'k}(\mathbf{x}_{t'+1}, t'+1) p_{t',j'} \Delta t \Delta t} \frac{d\{\mathbf{x}_t\} d\{\mathbf{p}_t\}}{(2\pi)^{DT}}$$

where we have used the shorthand $d\{\mathbf{x}_t\} d\{\mathbf{p}_t\} := d\mathbf{p}_0 d\mathbf{x}_1 d\mathbf{p}_1 \cdots d\mathbf{x}_{T-1} d\mathbf{p}_{T-1}$. This is our final answer, although it is more enlightening to write it in its schematic continuous-time form. We obtain

$$p(\mathbf{x}_0|\mathbf{x}_T) = \int \mathcal{D}[\mathbf{p}_t] \mathcal{D}[\mathbf{x}_t] \exp \left\{ \int_{\epsilon}^T -i\mathbf{p}_t \cdot [-\dot{\mathbf{x}}_t + \mathbf{f}(\mathbf{x}_t, t)] dt - \frac{1}{2} \int_{\epsilon}^T \int_{\epsilon}^T \mathbf{p}_t^T \mathbf{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') \mathbf{p}_{t'} dt dt' \right\}$$

where we have defined the state- and time-dependent $D \times D$ V-kernel $V_{ij}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t')$ via

$$V_{ij}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') := \sum_{k=1}^K G_{ik}(\mathbf{x}_t, t) G_{jk}(\mathbf{x}_{t'}, t'), \quad (49)$$

or equivalently via $\mathbf{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') := \mathbf{G}(\mathbf{x}_t, t) \mathbf{G}^T(\mathbf{x}_{t'}, t')$. Note that it is positive semidefinite.

D.3 AVERAGING LEARNED DISTRIBUTION OVER SAMPLE REALIZATIONS

What is the ‘typical’ distribution learned by an ensemble of diffusion models which differ only in the samples each used during training? In this subsection, we show that the net effect of averaging over sample realizations is to contribute a noise term to the PF-ODE. The path-integral representation we obtain is of the class we discussed in the previous subsection.

Suppose a diffusion model is associated with a parameterized score approximator $\hat{s}_\theta(\mathbf{x}, t)$. The distribution learned by the diffusion model is then

$$q(\mathbf{x}_0|\mathbf{x}_T; \boldsymbol{\theta}) = \int \mathcal{D}[\mathbf{p}_t] \mathcal{D}[\mathbf{x}_t] \exp \left\{ \int_\epsilon^T i\mathbf{p}_t \cdot [\dot{\mathbf{x}}_t + \beta_t \mathbf{x}_t + \mathbf{D}_t \hat{s}_\theta(\mathbf{x}_t, t)] dt \right\}, \quad (50)$$

where we have used the schematic form of the PF-ODE path-integral representation for clarity. (Moving to discrete time does not affect our arguments, but only makes notation more cumbersome.) Averaging over sample realizations is mathematically equivalent to computing the characteristic function of the score approximator. The sample-averaged q , $\mathbb{E}_\theta[q(\mathbf{x}_0|\mathbf{x}_T; \boldsymbol{\theta})] = [q(\mathbf{x}_0|\mathbf{x}_T)]$, is

$$[q(\mathbf{x}_0|\mathbf{x}_T)] = \int \mathcal{D}[\mathbf{p}_t] \mathcal{D}[\mathbf{x}_t] \exp \left\{ \int_\epsilon^T i\mathbf{p}_t \cdot [\dot{\mathbf{x}}_t + \beta_t \mathbf{x}_t] dt \right\} \mathbb{E}_\theta \left[e^{\int_\epsilon^T i\mathbf{p}_t^T \mathbf{D}_t \hat{s}_\theta(\mathbf{x}_t, t) dt} \right]. \quad (51)$$

Assuming the score approximator ensemble is well-behaved, its characteristic function can be written as a cumulant expansion. Here, we have

$$\begin{aligned} & \log \mathbb{E}_\theta \left[e^{\int_\epsilon^T i\mathbf{p}_t^T \mathbf{D}_t \hat{s}_\theta(\mathbf{x}_t, t) dt} \right] \\ &= \int_\epsilon^T i\mathbf{p}_t^T \mathbf{D}_t [\hat{s}_\theta(\mathbf{x}_t, t)] dt - \frac{1}{2} \int_\epsilon^T \int_\epsilon^T \mathbf{p}_t^T \mathbf{D}_t \text{Cov}_\theta [\hat{s}_\theta(\mathbf{x}_t, t), \hat{s}_\theta(\mathbf{x}_{t'}, t')] \mathbf{D}_{t'} \mathbf{p}_{t'} + \dots \end{aligned} \quad (52)$$

where the dots indicate higher-order cumulants and $[\hat{s}_\theta(\mathbf{x}_t, t)]$ indicates the ensemble-averaged score approximator. In this work, we neglect the higher-order terms. Often, they are suppressed by some factor (e.g., the number of model parameters divided by the number of samples).

We obtain dynamics of the class described in the previous subsection. Here, the $D \times D$ V-kernel is

$$V_{ij}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') := \sum_{a,b} D_{t,ia} \text{Cov}_\theta [\hat{s}_a(\mathbf{x}_t, t), \hat{s}_b(\mathbf{x}_{t'}, t')] D_{t',bj}, \quad (53)$$

or equivalently $\mathbf{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') := \mathbf{D}_t \text{Cov}_\theta [\hat{s}_\theta(\mathbf{x}_t, t), \hat{s}_\theta(\mathbf{x}_{t'}, t')] \mathbf{D}_{t'}$.

E NAIVE SCORE ESTIMATORS GENERALIZE: DETAILS

In this appendix, we show that integrating the PF-ODE using naive score estimates yields a specific kind of generalization (Prop. 4.1). Suppose that we are integrating the PF-ODE from some initial point \mathbf{x}_T using T first-order Euler updates (or some other integration scheme; the choice does not matter in the continuous-time limit), so that

$$\mathbf{x}_t = \mathbf{x}_{t+1} + (\beta_{t+1}\mathbf{x}_{t+1} + \mathbf{D}_{t+1}\mathbf{s}(\mathbf{x}_{t+1}, t + 1)) \Delta t. \quad (54)$$

But suppose that we do not use the score function directly in our updates, but at each time step construct a noisy version of it based on a sample $\mathbf{x}_{0t} \sim p(\mathbf{x}_0|\mathbf{x}, t)$. In particular, consider the naive score estimator

$$\hat{\mathbf{s}}(\mathbf{x}_t, t) := \mathbf{s}(\mathbf{x}_t, t) + \sqrt{\frac{\kappa}{\Delta t}} [\tilde{\mathbf{s}}(\mathbf{x}_t, t; \mathbf{x}_{0t}) - \mathbf{s}(\mathbf{x}_t, t)] \quad (55)$$

where $\kappa \geq 0$ is a constant that controls the estimator’s variance. Note that a new, independent sample $\mathbf{x}_{0t'}$ is drawn at each step t' . We are interested in studying the extent to which this scheme produces a distribution different from $p_{data}(\mathbf{x}_0)$.

Using the result from Appendix D, the typical learned distribution $[q(\mathbf{x}_0|\mathbf{x}_T)]$ is (approximately) characterized by the average and V-kernel of $\hat{\mathbf{s}}$. Since $\mathbb{E}_{\mathbf{x}_{0t}}[\tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_{0t})] = \mathbf{s}(\mathbf{x}, t)$, this score estimator is unbiased, i.e., $[\hat{\mathbf{s}}] = \mathbf{s}$. The V-kernel $\mathbf{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t')$ is

$$\mathbf{V} := \mathbf{D}_t \text{Cov}_{\theta}[\hat{\mathbf{s}}(\mathbf{x}_t, t), \hat{\mathbf{s}}(\mathbf{x}_{t'}, t')] \mathbf{D}_{t'} = \mathbf{D}_t \text{Cov}_{\theta}[\hat{\mathbf{s}}(\mathbf{x}_t, t), \hat{\mathbf{s}}(\mathbf{x}_t, t)] \mathbf{D}_t \delta(t - t') \Delta t \quad (56)$$

since samples generated at different time steps are independent of one another. Moreover,

$$\text{Cov}_{\theta}[\hat{\mathbf{s}}(\mathbf{x}_t, t)] = \frac{\kappa}{\Delta t} \text{Cov}_{\theta}[\tilde{\mathbf{s}}(\mathbf{x}_t, t)] \quad (57)$$

since the only random part of the estimator is the proxy score. Finally,

$$\begin{aligned} V_{ij}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') &= \kappa \sum_{a,b} D_{t,ia} \text{Cov}_{\theta}[\tilde{s}_a(\mathbf{x}_t, t), \tilde{s}_b(\mathbf{x}_t, t)] D_{t,bj} \delta(t - t') \\ &= \kappa \sum_{a,b} D_{t,ia} \left[S_{t,ab}^{-1} + \partial_{ab}^2 \log p(\mathbf{x}_t|t) \right] D_{t,bj} \delta(t - t'). \end{aligned} \quad (58)$$

Using $\mathbf{C}(\mathbf{x}, t)$ as shorthand for the proxy score covariance matrix, we equivalently have

$$\mathbf{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') = \kappa \mathbf{D}_t \mathbf{C}(\mathbf{x}_t, t) \mathbf{D}_t \delta(t - t'). \quad (59)$$

As a final technical note, note that the naive estimator must scale like $1/\sqrt{\Delta t}$ in order for the V-kernel to be nontrivial in the $\Delta t \rightarrow 0$ limit (and indeed, for the continuous-time limit to make sense). This is easiest to see in discrete time: since samples generated at different time steps k and ℓ are independent, the V-kernel picks up a factor $\delta_{k\ell}$, which equals one when $k = \ell$ and is zero otherwise. In continuous time, this looks like $\delta(t - t') \Delta t$, *not* $\delta(t - t')$ (Eq. 56). This problematic Δt factor can be canceled by a corresponding $1/(\Delta t)$ factor in the estimator covariance, which motivates making the estimator scale like $1/\sqrt{\Delta t}$.

F LINEAR SCORE ESTIMATOR: DETAILS

In this appendix, we compute the sample-realization-averaged distribution learned by a linear score estimator (Prop. 5.1). Whether it generalizes or not depends strongly on whether the number of features F scales with the number of samples P used during training. First, we must compute the optimum of the DSM objective for a linear model. Then we will determine the average and V-kernel of the optimal linear score estimator.

F.1 DEFINITION OF LINEAR SCORE MODEL

Consider a linear score estimator

$$\hat{\mathbf{s}}_{\theta}(\mathbf{x}, t) = \mathbf{w}_0 + \mathbf{W}\phi(\mathbf{x}, t) \quad \hat{s}_i(\mathbf{x}, t) = w_{0i} + \sum_{j=1}^F W_{ij}\phi_j(\mathbf{x}, t), \quad (60)$$

where the feature maps $\phi = (\phi_1, \dots, \phi_F)^T$ are linearly independent, smooth functions from $\mathbb{R}^D \times [0, T]$ to \mathbb{R} that are square-integrable with respect to the measure $\lambda_t p(\mathbf{x}, t)$ for all t . The parameters to be estimated are $\theta := \{\mathbf{w}_0, \mathbf{W}\}$, with $\mathbf{w}_0 \in \mathbb{R}^D$ and $\mathbf{W} \in \mathbb{R}^{D \times F}$.

F.2 OPTIMUM OF DSM OBJECTIVE FOR LINEAR SCORE MODEL

For this estimator, the DSM objective reads

$$J_1(\theta) = \int \frac{\lambda_t}{2} \|\mathbf{w}_0 + \mathbf{W}\phi(\mathbf{x}, t) - \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)\|_2^2 p(\mathbf{x}|\mathbf{x}_0, t) p_{data}(\mathbf{x}_0) p(t) d\mathbf{x} d\mathbf{x}_0 dt. \quad (61)$$

Note,

$$\frac{\partial \hat{s}_i}{\partial w_{0a}} = \delta_{ia} \quad \frac{\partial \hat{s}_i}{\partial W_{ab}} = \delta_{ia} \phi_b. \quad (62)$$

Using these to take the gradient of the DSM objective, we have

$$\begin{aligned} \frac{\partial J_1}{\partial w_{0a}} &= \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \left\{ \lambda_t \left[w_{0a} + \sum_{j=1}^F W_{aj} \phi_j(\mathbf{x}, t) - \tilde{s}_a(\mathbf{x}, t; \mathbf{x}_0) \right] \right\} \\ \frac{\partial J_1}{\partial W_{ab}} &= \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \left\{ \lambda_t \left[w_{0a} + \sum_{j=1}^F W_{aj} \phi_j(\mathbf{x}, t) - \tilde{s}_a(\mathbf{x}, t; \mathbf{x}_0) \right] \phi_b(\mathbf{x}, t) \right\}. \end{aligned} \quad (63)$$

Setting these equal to zero, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \} w_{0a} + \sum_{j=1}^F W_{aj} \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \phi_j(\mathbf{x}, t) \} &= \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \tilde{s}_a(\mathbf{x}, t; \mathbf{x}_0) \} \\ \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \phi_b(\mathbf{x}, t) \} w_{0a} + \sum_{j=1}^F W_{aj} \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \phi_j(\mathbf{x}, t) \phi_b(\mathbf{x}, t) \} &= \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \tilde{s}_a(\mathbf{x}, t; \mathbf{x}_0) \phi_b(\mathbf{x}, t) \}. \end{aligned}$$

The first row tells us that

$$w_{0a} = \frac{1}{\mathbb{E}_t[\lambda_t]} \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \tilde{s}_a(\mathbf{x}, t; \mathbf{x}_0) \} - \frac{1}{\mathbb{E}_t[\lambda_t]} \sum_{j=1}^F W_{aj} \mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \phi_j(\mathbf{x}, t) \}, \quad (64)$$

or equivalently that the optimal bias term satisfies $\mathbf{w}_0^* = \langle \tilde{\mathbf{s}} \rangle - \mathbf{W}^* \langle \phi \rangle$, where we have used $\langle \dots \rangle$ to denote averages with respect to $\lambda_t p(\mathbf{x}, \mathbf{x}_0, t) / \mathbb{E}[\lambda_t]$, and where we have defined the vectors

$$\begin{aligned} \langle \tilde{\mathbf{s}} \rangle &:= \frac{\mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} [\lambda_t \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)]}{\mathbb{E}_t[\lambda_t]} = \frac{1}{\mathbb{E}_t[\lambda_t]} \int \lambda_t \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0) p(\mathbf{x}, \mathbf{x}_0, t) d\mathbf{x} d\mathbf{x}_0 dt \\ \langle \phi \rangle &:= \frac{\mathbb{E}_{\mathbf{x}, t} [\lambda_t \phi(\mathbf{x}, t)]}{\mathbb{E}_t[\lambda_t]} = \frac{1}{\mathbb{E}_t[\lambda_t]} \int \lambda_t \phi(\mathbf{x}, t) p(\mathbf{x}, t) d\mathbf{x} dt. \end{aligned} \quad (65)$$

Using the first row result, the second row can be written as

$$\langle \phi_b \rangle \left[\langle \tilde{s}_a \rangle - \sum_{j=1}^F W_{aj} \langle \phi_j \rangle \right] + \sum_{j=1}^F W_{aj} \frac{\mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \phi_j(\mathbf{x}, t) \phi_b(\mathbf{x}, t) \}}{\mathbb{E}_t[\lambda_t]} = \frac{\mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \tilde{s}_a(\mathbf{x}, t; \mathbf{x}_0) \phi_b(\mathbf{x}, t) \}}{\mathbb{E}_t[\lambda_t]}$$

and hence the second row can be written in terms of matrices

$$\begin{aligned} \Sigma_\phi &:= \frac{\mathbb{E}_{\mathbf{x}, t} \{ \lambda_t [\phi(\mathbf{x}, t) - \langle \phi \rangle] [\phi(\mathbf{x}, t) - \langle \phi \rangle]^T \}}{\mathbb{E}_t[\lambda_t]} \\ &= \frac{1}{\mathbb{E}_t[\lambda_t]} \int \lambda_t [\phi(\mathbf{x}, t) - \langle \phi \rangle] [\phi(\mathbf{x}, t) - \langle \phi \rangle]^T p(\mathbf{x}, t) d\mathbf{x} dt \\ \mathbf{J} &:= - \frac{\mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t [\phi(\mathbf{x}, t) - \langle \phi \rangle] [\tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0) - \langle \tilde{\mathbf{s}} \rangle]^T \}}{\mathbb{E}_t[\lambda_t]} \\ &= - \frac{1}{\mathbb{E}_t[\lambda_t]} \int \lambda_t [\phi(\mathbf{x}, t) - \langle \phi \rangle] [\tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0) - \langle \tilde{\mathbf{s}} \rangle]^T p(\mathbf{x}, \mathbf{x}_0, t) d\mathbf{x} d\mathbf{x}_0 dt. \end{aligned} \quad (66)$$

In particular,

$$\mathbf{W}^* \Sigma_\phi = -\mathbf{J}^T \implies \mathbf{W}^* = -\mathbf{J}^T \Sigma_\phi^{-1}, \quad (67)$$

where we have assumed that Σ_ϕ is invertible. This ought to be true, since the feature maps are independent and $p(\mathbf{x}|t)$ is a smooth distribution supported on all of \mathbb{R}^D (especially since we are technically only considering t as small as ϵ , the nonzero lower bound, for regularization purposes).

The optimal score is

$$\hat{s}_*(\mathbf{x}, t) = \mathbf{w}_0^* + \mathbf{W}^* \phi(\mathbf{x}, t) = \mathbf{J}^T \Sigma_\phi^{-1} [\langle \phi \rangle - \phi(\mathbf{x}, t)] + \langle \tilde{\mathbf{s}} \rangle. \quad (68)$$

As a side comment, omitting the bias term just removes the mean corrections from the definitions of \mathbf{J} and Σ_ϕ , as well as the $\langle \phi \rangle$ and $\langle \tilde{\mathbf{s}} \rangle$ offsets. Without it, the optimal score is $\hat{s}_*(\mathbf{x}, t) = \mathbf{W}^* \phi(\mathbf{x}, t) = -\mathbf{J}^T \Sigma_\phi^{-1} \phi(\mathbf{x}, t)$, where \mathbf{J} and Σ_ϕ are instead defined to be

$$\begin{aligned} \Sigma_\phi &:= \frac{\mathbb{E}_{\mathbf{x}, t} \{ \lambda_t \phi(\mathbf{x}, t) \phi(\mathbf{x}, t)^T \}}{\mathbb{E}_t[\lambda_t]} \\ \mathbf{J} &:= - \frac{\mathbb{E}_{\mathbf{x}, \mathbf{x}_0, t} \{ \lambda_t \phi(\mathbf{x}, t) \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)^T \}}{\mathbb{E}_t[\lambda_t]}. \end{aligned} \quad (69)$$

In the rest of this appendix, we will assume that the bias term is present.

F.3 OPTIMUM OF DSM OBJECTIVE GIVEN A FINITE NUMBER OF SAMPLES

Assume we have access to $P \gg 1$ samples $\mathbf{x}^{(n)}, \mathbf{x}_0^{(n)}, t^{(n)} \sim p(\mathbf{x}, \mathbf{x}_0, t)$, and that we estimate the parameters of the linear score model using naive sample mean estimators

$$\begin{aligned} \bar{\lambda}_t &:= \frac{1}{P} \sum_n \lambda^{(n)} \\ \hat{\mathbf{b}} &:= \frac{1}{\bar{\lambda}_t} \frac{1}{P} \sum_n \lambda^{(n)} \tilde{\mathbf{s}}(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}_0^{(n)}) \\ \hat{\boldsymbol{\mu}}_\phi &:= \frac{1}{\bar{\lambda}_t} \frac{1}{P} \sum_n \lambda^{(n)} \phi(\mathbf{x}^{(n)}, t^{(n)}) \\ \hat{\Sigma}_\phi &:= \frac{1}{\bar{\lambda}_t} \frac{1}{P} \sum_n \lambda^{(n)} \left[\phi(\mathbf{x}^{(n)}, t^{(n)}) - \hat{\boldsymbol{\mu}}_\phi \right] \left[\phi(\mathbf{x}^{(n)}, t^{(n)}) - \hat{\boldsymbol{\mu}}_\phi \right]^T \\ \hat{\mathbf{J}} &:= - \frac{1}{\bar{\lambda}_t} \frac{1}{P} \sum_n \lambda^{(n)} \left[\phi(\mathbf{x}^{(n)}, t^{(n)}) - \hat{\boldsymbol{\mu}}_\phi \right] \left[\tilde{\mathbf{s}}(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}_0^{(n)}) - \hat{\mathbf{b}} \right]^T \end{aligned} \quad (70)$$

where we have used $\lambda^{(n)}$ as a slightly less cumbersome shorthand for $\lambda_{t^{(n)}}$. We will not worry about using Bessel's correction in the covariance estimators, and we will see below that $\hat{\mathbf{s}}$ is actually unbiased for finite P even if the covariance estimators are not. Our learned score estimator is then

$$\hat{s}_\theta(\mathbf{x}, t) = \hat{\mathbf{J}}^T \hat{\Sigma}_\phi^{-1} [\hat{\boldsymbol{\mu}}_\phi - \phi(\mathbf{x}, t)] + \hat{\mathbf{b}}. \quad (71)$$

Note that if the number of samples P is less than F , $\hat{\Sigma}_\phi$ is not invertible; in this case, one can either use the Moore-Penrose pseudoinverse, or explicitly include a weight regularization term in the objective, i.e., add to J_1 a term of the form

$$J_{reg} := \frac{\xi}{2} \left[\sum_i w_{0i}^2 + \sum_{i,j} W_{ij}^2 \right] \quad (72)$$

where the parameter $\xi \geq 0$ controls the importance of this term. Including this term changes the score estimator (Eq. 71) by modifying the inverse that appears:

$$\hat{\Sigma}_\phi^{-1} \rightarrow \left[\hat{\Sigma}_\phi + \xi \mathbf{I}_F \right]^{-1}. \quad (73)$$

In what follows, if one wants results in the case that such a regularization term is present, note that this replacement of the inverse of the empirical covariance matrix is the only necessary change.

F.4 LINEAR SCORE MODEL ESTIMATOR IS UNBIASED

We are primarily interested in variance due to \mathbf{x}_0 (for reasons that will become clear), so we will consider an ensemble of systems for which the $\mathbf{x}^{(n)}$ and $t^{(n)}$ sample draws are the same, but the $\mathbf{x}_0^{(n)}$ draws are different. Our estimator depends linearly on $\tilde{\mathbf{s}}$, the quantity through which it depends on the \mathbf{x}_0 samples. In particular,

$$\begin{aligned} \hat{\mathbf{J}}^T \hat{\Sigma}_\phi^{-1} [\hat{\boldsymbol{\mu}}_\phi - \phi(\mathbf{x}, t)] &= \frac{1}{P} \sum_n \frac{\lambda^{(n)}}{\lambda_t} \left[\tilde{\mathbf{s}}(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}_0^{(n)}) - \hat{\mathbf{b}} \right] \left[\phi(\mathbf{x}^{(n)}, t^{(n)}) - \hat{\boldsymbol{\mu}}_\phi \right]^T \hat{\Sigma}_\phi^{-1} [\phi(\mathbf{x}, t) - \hat{\boldsymbol{\mu}}_\phi] \\ &= \frac{1}{P} \sum_n \frac{\lambda^{(n)}}{\lambda_t} \tilde{\mathbf{s}}(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}_0^{(n)}) \left[\phi(\mathbf{x}^{(n)}, t^{(n)}) - \hat{\boldsymbol{\mu}}_\phi \right]^T \hat{\Sigma}_\phi^{-1} [\phi(\mathbf{x}, t) - \hat{\boldsymbol{\mu}}_\phi], \end{aligned}$$

so

$$\hat{\mathbf{s}}_\theta(\mathbf{x}, t) = \frac{1}{P} \sum_n \frac{\lambda^{(n)}}{\lambda_t} Q(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}, t) \tilde{\mathbf{s}}(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}_0^{(n)}) \quad (74)$$

where we have defined the kernel function

$$Q(\mathbf{x}, t; \mathbf{x}', t') := 1 + [\phi(\mathbf{x}, t) - \hat{\boldsymbol{\mu}}]^T \hat{\Sigma}_\phi^{-1} [\phi(\mathbf{x}', t') - \hat{\boldsymbol{\mu}}]. \quad (75)$$

To see that this estimator is unbiased (when the model is sufficiently expressive), suppose the true score has the form of our linear estimator, i.e.,

$$\mathbf{s}(\mathbf{x}, t) = \mathbf{w}_0^* + \mathbf{W}^* \phi(\mathbf{x}, t) = \mathbf{W}^* [\phi(\mathbf{x}, t) - \langle \phi \rangle], \quad (76)$$

where we have used the fact that $\mathbb{E}_\mathbf{x}[\mathbf{s}] = \langle \mathbf{s} \rangle = \mathbf{0}$. Next, note that

$$\frac{1}{P} \sum_n \frac{\lambda^{(n)}}{\lambda_t} Q(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}, t) = 1. \quad (77)$$

Averaging our estimator over \mathbf{x}_0 sample draws yields

$$\begin{aligned} \mathbb{E}[\hat{\mathbf{s}}_\theta(\mathbf{x}, t)] &= \frac{1}{P} \sum_n \frac{\lambda^{(n)}}{\lambda_t} Q(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}, t) \mathbf{W}^* [\phi(\mathbf{x}^{(n)}, t^{(n)}) - \hat{\boldsymbol{\mu}} + \hat{\boldsymbol{\mu}} - \langle \phi \rangle] \\ &= \frac{1}{P} \sum_n \frac{\lambda^{(n)}}{\lambda_t} Q(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}, t) \mathbf{W}^* [\phi(\mathbf{x}^{(n)}, t^{(n)}) - \hat{\boldsymbol{\mu}}] + \mathbf{W}^* (\hat{\boldsymbol{\mu}} - \langle \phi \rangle) \\ &= \frac{1}{P} \sum_n \frac{\lambda^{(n)}}{\lambda_t} \mathbf{W}^* [\phi(\mathbf{x}^{(n)}, t^{(n)}) - \hat{\boldsymbol{\mu}}] \left[\phi(\mathbf{x}^{(n)}, t^{(n)}) - \hat{\boldsymbol{\mu}} \right]^T \hat{\Sigma}_\phi^{-1} (\phi(\mathbf{x}, t) - \hat{\boldsymbol{\mu}}) + \mathbf{W}^* (\hat{\boldsymbol{\mu}} - \langle \phi \rangle) \\ &= \mathbf{W}^* (\phi(\mathbf{x}, t) - \hat{\boldsymbol{\mu}}) + \mathbf{W}^* (\hat{\boldsymbol{\mu}} - \langle \phi \rangle) \\ &= \mathbf{W}^* (\phi(\mathbf{x}, t) - \langle \phi \rangle), \end{aligned}$$

i.e., it is unbiased. What is worth emphasizing is that this is *exactly* true, and does not require taking any kind of large P limit. In other words, as long as P is large enough that $\hat{\Sigma}_\phi$ is invertible, one recovers the true weights \mathbf{w}_0^* and \mathbf{W}^* , independent of the \mathbf{x} and t sample draws. This is why variance due to \mathbf{x}_0 sample draws matters, and variance due to the other draws does not, at least for this linear model.

F.5 COMPUTING THE V-KERNEL OF THE LINEAR SCORE MODEL

Computing the V-kernel amounts to computing the covariance of the score model with respect to \mathbf{x}_0 sample realizations. In the previous section, we computed the mean of our estimator; the covariance calculation will be fairly similar. Note that

$$\begin{aligned} \text{Cov}[\hat{\mathbf{s}}(\mathbf{z}), \hat{\mathbf{s}}(\mathbf{z}')] &= \frac{1}{P^2} \sum_{n,m} \frac{\lambda^{(n)}}{\bar{\lambda}_t} \frac{\lambda^{(m)}}{\bar{\lambda}_t} Q(\mathbf{z}^{(n)}; \mathbf{z}) Q(\mathbf{z}^{(m)}; \mathbf{z}') \text{Cov}[\tilde{\mathbf{s}}(\mathbf{z}^{(n)}; \mathbf{x}_0^{(n)}), \tilde{\mathbf{s}}(\mathbf{z}^{(m)}; \mathbf{x}_0^{(m)})] \\ &= \frac{1}{P^2} \sum_n \left(\frac{\lambda^{(n)}}{\bar{\lambda}_t} \right)^2 Q(\mathbf{z}^{(n)}; \mathbf{z}) Q(\mathbf{z}^{(n)}; \mathbf{z}') \text{Cov}[\tilde{\mathbf{s}}(\mathbf{z}^{(n)}; \mathbf{x}_0^{(n)})], \end{aligned}$$

where we have used \mathbf{z} as shorthand for $\{\mathbf{x}, t\}$, and the fact that \mathbf{x}_0 sample draws are independent of one another. Now we will invoke the central limit theorem. Using $\mathbf{C}(\mathbf{z}) := \text{Cov}[\tilde{\mathbf{s}}(\mathbf{z}; \mathbf{x}_0)]$ as shorthand, when P is very large, to leading order in $1/P$ we have

$$\begin{aligned} \text{Cov}[\hat{\mathbf{s}}(\mathbf{z}), \hat{\mathbf{s}}(\mathbf{z}')] &\approx \frac{1}{P} \int \left(\frac{\lambda_t}{\bar{\lambda}_t} \right)^2 Q(\mathbf{z}''; \mathbf{z}) Q(\mathbf{z}''; \mathbf{z}') \mathbf{C}(\mathbf{z}'') p(\mathbf{z}'') d\mathbf{z}'' \\ &\approx \frac{1}{P} \int \left(\frac{\lambda_t}{\mathbb{E}[\lambda_t]} \right)^2 Q(\mathbf{z}''; \mathbf{z}) Q(\mathbf{z}''; \mathbf{z}') \mathbf{C}(\mathbf{z}'') p(\mathbf{z}'') d\mathbf{z}'' \end{aligned} \quad (78)$$

where we replace the estimates $\hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\Sigma}}_\phi$ that appear in the kernel function with the true quantities, i.e., we redefine Q to be

$$Q(\mathbf{x}'', t''; \mathbf{x}, t) := 1 + [\boldsymbol{\phi}(\mathbf{x}'', t'') - \langle \boldsymbol{\phi} \rangle]^T \boldsymbol{\Sigma}_\phi^{-1} [\boldsymbol{\phi}(\mathbf{x}, t) - \langle \boldsymbol{\phi} \rangle]. \quad (79)$$

If the number of features F does *not* scale with the number of samples P , then we are done: in the $P \rightarrow \infty$ limit, the score estimator covariance, and hence the V-kernel, approach zero. Alternatively, if the number of features F *does* scale with the number of samples P , a nontrivial result is possible.

The second term of Q , a quadratic form involving the model's feature maps, is the only place in Eq. 78 one can get nontrivial scaling with F . Motivated by this observation, define the feature kernel

$$k(\mathbf{z}; \mathbf{z}') := \frac{1}{\sqrt{F}} [\boldsymbol{\phi}(\mathbf{z}) - \langle \boldsymbol{\phi} \rangle]^T \boldsymbol{\Sigma}_\phi^{-1} [\boldsymbol{\phi}(\mathbf{z}') - \langle \boldsymbol{\phi} \rangle]. \quad (80)$$

Provided that the limit exists and is finite, in the $P \rightarrow \infty$ limit (where F may scale with P), the asymptotic V-kernel is then

$$\mathbf{V}(\mathbf{z}; \mathbf{z}') = \lim_{P \rightarrow \infty} \frac{F}{P} \mathbf{D}_t \mathbb{E}_{\mathbf{z}''} \left\{ \frac{\lambda_{t''}^2}{\mathbb{E}_t[\lambda_t]^2} k(\mathbf{z}; \mathbf{z}'') \mathbf{C}(\mathbf{z}'') k(\mathbf{z}''; \mathbf{z}') \right\} \mathbf{D}_{t'}. \quad (81)$$

Note also that, in the large P limit, the V-kernel *also* does not depend on the \mathbf{x} and t sample draws.

G NEURAL NETWORK SCORE ESTIMATOR IN NTK REGIME: DETAILS

In this appendix, we prove Prop. 5.2, which means computing the V-kernel of a fully-connected, infinite-width neural network in the ‘lazy’ learning (Chizat et al., 2019) regime. Although we focus on an extremely specific type of network here, note that our argument can be straightforwardly adapted to compute the V-kernel of other architectures with NTK limits, like convolutional neural networks (Arora et al., 2019).

G.1 DEFINITION OF NEURAL NETWORK MODEL

Consider a neural network score function approximator $\hat{s}_\theta(\mathbf{x}, t)$ trained on the DSM objective (Eq. 4). As elsewhere, we may use \mathbf{z} as shorthand for $\{\mathbf{x}, t\}$. For concreteness, assume that the network is fully-connected, has $L \geq 1$ layers and N_* trainable parameters, and that each hidden layer has N neurons and an identical pointwise nonlinearity G :

$$\begin{aligned} a_i^{(0)}(\mathbf{z}) &:= \psi_i(\mathbf{z}) \\ a_i^{(\ell+1)}(\mathbf{z}) &:= G\left(\frac{1}{\sqrt{N}} \sum_j W_{ij}^{(\ell+1)} a_j^{(\ell)}(\mathbf{z})\right) \quad \ell = 0, \dots, L-2 \\ a_i^{(L)}(\mathbf{z}) &:= \frac{1}{\sqrt{N}} \sum_j W_{ij}^{(L)} a_j^{(L-1)}(\mathbf{z}) \quad \hat{s}_\theta(\mathbf{z}) := \mathbf{a}^{(L)}(\mathbf{z}). \end{aligned} \quad (82)$$

The (non-trainable) initial feature maps $\boldsymbol{\psi} := (\psi_1, \dots, \psi_{N_0})^T$ account for various preconditioning-related choices. For example, in practice, diffusion models receive time/noise as input only through some time/noise embedding (Ho et al., 2020; Song et al., 2021; Karras et al., 2022).

Although characterizing the gradient descent dynamics of \hat{s} may be difficult in general, if the initial network weights are sampled i.i.d. from a standard normal (i.e., $W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1)$ for all i, j , and ℓ), as N is taken to infinity the network output becomes independent of the precise values of the initial weights. Moreover, the network’s output throughout training can be written in terms of a kernel function—the so-called NTK—defined by

$$K^{cc'}(\mathbf{z}, \mathbf{z}') := \sum_i \mathbb{E}_\theta \left\{ \frac{\partial \hat{s}_c(\mathbf{z})}{\partial \theta_i} \frac{\partial \hat{s}_{c'}(\mathbf{z}')}{\partial \theta_i} \right\} \quad (83)$$

where c and c' index different network outputs. In the infinite-width ($N \rightarrow \infty$) limit, $K^{cc'}(\mathbf{z}, \mathbf{z}') = \delta_{cc'} K(\mathbf{z}, \mathbf{z}')$, i.e., the off-diagonal kernels are identically zero and all kernels along the diagonal are the same (Shan & Bordelon, 2022).

G.2 LEARNED SCORE AFTER FULL-BATCH GRADIENT DESCENT

Computing the learned score. For simplicity, we assume that our neural network model is trained via full-batch gradient descent on P samples from $p(\mathbf{x}, \mathbf{x}_0, t)$. Although this assumption does not reflect standard practice (Song et al., 2021; Karras et al., 2022), it makes our computation substantially easier. If we let the dimensionless parameter τ denote training time, the output evolves via

$$\frac{d}{d\tau} \hat{s}(\mathbf{x}', t') = \mathbb{E}_{\mathbf{x}, t, \mathbf{x}_0} \left\{ \frac{\lambda_t}{\mathbb{E}_t[\lambda_t]} \frac{\partial \hat{s}(\mathbf{x}', t')}{\partial \boldsymbol{\theta}} \frac{\partial \hat{s}(\mathbf{x}, t)}{\partial \boldsymbol{\theta}} [\tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0) - \hat{s}(\mathbf{x}, t)] \right\}. \quad (84)$$

In the infinite-width limit, we can replace the outer product that appears with the NTK:

$$\frac{d}{d\tau} \hat{s}(\mathbf{x}', t') = \mathbb{E}_{\mathbf{x}, t, \mathbf{x}_0} \left\{ \frac{\lambda_t}{\mathbb{E}_t[\lambda_t]} K(\mathbf{x}', t'; \mathbf{x}, t) [\tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0) - \hat{s}(\mathbf{x}, t)] \right\}. \quad (85)$$

Define the Gram matrix $\mathbf{K} \in \mathbb{R}^{P \times P}$, the time-weighting matrix $\mathbf{\Lambda}_T \in \mathbb{R}^{P \times P}$, the target matrix $\tilde{\mathbf{S}} \in \mathbb{R}^{P \times D}$, and the output matrix $\hat{\mathbf{S}} \in \mathbb{R}^{P \times D}$ via

$$\begin{aligned} K_{ab} &:= K(\mathbf{x}^{(a)}, t^{(a)}; \mathbf{x}^{(b)}, t^{(b)}) \\ \Lambda_{T,ab} &:= \delta_{ab} \frac{\lambda_{t^{(a)}}}{\mathbb{E}_t[\lambda_t]} \\ \tilde{S}_{ai} &:= \tilde{s}_i(\mathbf{x}^{(a)}, t^{(a)}; \mathbf{x}_0^{(a)}) \\ \hat{S}_{ai} &:= \hat{s}_i(\mathbf{x}^{(a)}, t^{(a)}) . \end{aligned} \quad (86)$$

Eq. 85 implies that

$$\frac{d}{d\tau} \hat{\mathbf{S}} = \frac{1}{P} \mathbf{K} \mathbf{\Lambda}_T (\tilde{\mathbf{S}} - \hat{\mathbf{S}}) . \quad (87)$$

Hence, after training, the network's output on the set of samples is given by

$$\hat{\mathbf{S}} = e^{-\mathbf{K} \mathbf{\Lambda}_T \tau / P} \hat{\mathbf{S}}_0 + (\mathbf{I} - e^{-\mathbf{K} \mathbf{\Lambda}_T \tau / P}) \tilde{\mathbf{S}} \quad (88)$$

where τ is the total training 'time' and $\hat{\mathbf{S}}_0$ is the $P \times D$ matrix containing the network's initial output on the samples. Let $\mathbf{k}(\mathbf{x}, t)$ denote the P -dimensional vector whose i -th component is $K(\mathbf{x}^{(i)}, t^{(i)}; \mathbf{x}, t)$. The network's output given other inputs evolves according to the ODE

$$\frac{d}{d\tau} \hat{\mathbf{s}}(\mathbf{x}, t)^T = \frac{1}{P} \mathbf{k}(\mathbf{x}, t)^T \mathbf{\Lambda}_T (\tilde{\mathbf{S}} - \hat{\mathbf{S}}) , \quad (89)$$

whose solution is

$$\hat{\mathbf{s}}(\mathbf{x}, t)^T = \hat{\mathbf{s}}_0(\mathbf{x}, t)^T + \mathbf{k}(\mathbf{x}, t)^T \mathbf{K}^{-1} (\mathbf{I} - e^{-\mathbf{K} \mathbf{\Lambda}_T \tau / P}) (\tilde{\mathbf{S}} - \hat{\mathbf{S}}_0) \quad (90)$$

where $\hat{\mathbf{s}}_0(\mathbf{x}, t)$ is the network's initial output given a $\{\mathbf{x}, t\}$ input. If the Gram matrix \mathbf{K} is rank-deficient, we must use its Moore-Penrose pseudoinverse. Alternatively, one can avoid this issue by including a weight regularization term in the objective.

Expressing the learned score in terms of eigenfunctions. We will find it useful to consider a Mercer decomposition of K with respect to the measure $\lambda_t p(\mathbf{x}, t) / \mathbb{E}_t[\lambda_t]$, so that K can be written

$$K(\mathbf{x}, t; \mathbf{x}', t') = \sum_k \lambda_k \phi_k(\mathbf{x}, t) \phi_k(\mathbf{x}', t') \quad (91)$$

where the features are orthonormal and complete, i.e.,

$$\begin{aligned} \int \frac{\lambda_t}{\mathbb{E}_t[\lambda_t]} \phi_k(\mathbf{x}, t) \phi_{k'}(\mathbf{x}, t) p(\mathbf{x}, t) d\mathbf{x} dt &= \delta_{k,k'} \\ \sum_k \frac{\lambda_t}{\mathbb{E}_t[\lambda_t]} p(\mathbf{x}, t) \phi_k(\mathbf{x}, t) \phi_k(\mathbf{x}', t') &= \delta(\mathbf{x} - \mathbf{x}') \delta(t - t') . \end{aligned} \quad (92)$$

If we assume \mathbf{K} has rank F not necessarily equal to P , we can write Eq. 90 in terms of the eigenfunctions associated with the Mercer decomposition by defining the $P \times F$ matrix $\mathbf{\Phi}$ with

$$\Phi_{ak} := \phi_k(\mathbf{x}^{(a)}, t^{(a)}) \quad (93)$$

and noting that $\mathbf{K} = \mathbf{\Phi} \mathbf{\Lambda} \mathbf{\Phi}^T$, where $\mathbf{\Lambda}$ is the $F \times F$ diagonal matrix of associated eigenvalues. It is useful to observe that

$$\delta_{kk'} = \mathbb{E}_{\mathbf{x}, t} \left\{ \frac{\lambda_t}{\mathbb{E}_t[\lambda_t]} \phi_k(\mathbf{x}, t) \phi_{k'}(\mathbf{x}, t) \right\} = \frac{1}{P} \sum_n \frac{\lambda^{(n)}}{\mathbb{E}_t[\lambda_t]} \phi_k(\mathbf{x}^{(n)}, t^{(n)}) \phi_{k'}(\mathbf{x}^{(n)}, t^{(n)}) + \mathcal{O}(1/\sqrt{P}) ,$$

which implies $\mathbf{I}_F = \frac{\mathbf{\Phi}^T \mathbf{\Lambda}_T \mathbf{\Phi}}{P}$ to leading order. Similarly, the completeness relation becomes

$$\mathbf{I}_P \approx \frac{\mathbf{\Phi} \mathbf{\Phi}^T \mathbf{\Lambda}_T}{P} = \frac{\mathbf{\Lambda}_T \mathbf{\Phi} \mathbf{\Phi}^T}{P} \quad (94)$$

to leading order. Using these identities, we can rewrite Eq. 90 as

$$\begin{aligned}\hat{\mathbf{s}}(\mathbf{x}, t)^T &= \hat{\mathbf{s}}_0(\mathbf{x}, t)^T + [\phi(\mathbf{x}, t)^T \Lambda \Phi^T] \left[\frac{\Lambda_T \Phi \Lambda^{-1} \Phi^T \Lambda_T}{P^2} \right] \left[\frac{\Phi}{P} (\mathbf{I} - e^{-\Lambda \tau}) \Phi^T \Lambda_T \right] (\tilde{\mathbf{S}} - \hat{\mathbf{S}}_0) \\ &= \hat{\mathbf{s}}_0(\mathbf{x}, t)^T + \frac{1}{P} \phi(\mathbf{x}, t)^T (\mathbf{I} - e^{-\Lambda \tau}) \Phi^T \Lambda_T (\tilde{\mathbf{S}} - \hat{\mathbf{S}}_0).\end{aligned}\quad (95)$$

Equivalently,

$$\hat{\mathbf{s}}(\mathbf{x}, t) = \hat{\mathbf{s}}_0(\mathbf{x}, t) + \frac{1}{P} (\tilde{\mathbf{S}} - \hat{\mathbf{S}}_0)^T \Lambda_T \Phi (\mathbf{I} - e^{-\Lambda \tau}) \phi(\mathbf{x}, t). \quad (96)$$

Let \mathbf{S} denote the $P \times D$ matrix whose entries are the true score evaluated on the set of input samples $\{(\mathbf{x}^{(a)}, t^{(a)})\}$. When averaged over \mathbf{x}_0 sample realizations, our estimator is

$$\mathbb{E}[\hat{\mathbf{s}}(\mathbf{x}, t)] = \hat{\mathbf{s}}_0(\mathbf{x}, t) + \frac{1}{P} (\mathbf{S} - \hat{\mathbf{S}}_0)^T \Lambda_T \Phi (\mathbf{I} - e^{-\Lambda \tau}) \phi(\mathbf{x}, t) \quad (97)$$

which implies

$$\hat{\mathbf{s}}(\mathbf{x}, t) - \mathbb{E}[\hat{\mathbf{s}}(\mathbf{x}, t)] = \frac{1}{P} (\tilde{\mathbf{S}} - \mathbf{S})^T \Lambda_T \Phi (\mathbf{I} - e^{-\Lambda \tau}) \phi(\mathbf{x}, t). \quad (98)$$

To make things slightly easier, define the feature kernel¹

$$k(\mathbf{x}, t; \mathbf{x}', t') := \frac{1}{\sqrt{F}} \phi(\mathbf{x}, t)^T (\mathbf{I} - e^{-\Lambda \tau}) \phi(\mathbf{x}', t') = \frac{1}{\sqrt{F}} \sum_{k=1}^F \phi_k(\mathbf{x}, t) (1 - e^{-\lambda_k \tau}) \phi_k(\mathbf{x}', t'). \quad (99)$$

In terms of this kernel, we can write

$$\hat{\mathbf{s}}(\mathbf{x}, t) - \mathbb{E}[\hat{\mathbf{s}}(\mathbf{x}, t)] = \frac{\sqrt{F}}{P} \sum_n \frac{\lambda^{(n)}}{\mathbb{E}_t[\lambda_t]} \left[\tilde{\mathbf{s}}(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}_0^{(n)}) - \mathbf{s}(\mathbf{x}^{(n)}, t^{(n)}) \right] k(\mathbf{x}^{(n)}, t^{(n)}; \mathbf{x}, t). \quad (100)$$

We will use this result in the next subsection to compute the V-kernel of this model.

G.3 COMPUTING THE V-KERNEL OF THE NTK MODEL

The covariance of the learned score estimator with respect to \mathbf{x}_0 sample realizations is

$$\begin{aligned}\text{Cov}[\hat{\mathbf{s}}_\theta(\mathbf{z}), \hat{\mathbf{s}}_\theta(\mathbf{z}')] &= \frac{F}{P^2} \sum_{n,m} \frac{\lambda^{(n)}}{\mathbb{E}_t[\lambda_t]} \frac{\lambda^{(m)}}{\mathbb{E}_t[\lambda_t]} \text{Cov} \left[\tilde{\mathbf{s}}(\mathbf{z}^{(n)}; \mathbf{x}_0^{(n)}), \tilde{\mathbf{s}}(\mathbf{z}^{(m)}; \mathbf{x}_0^{(m)}) \right] k(\mathbf{z}^{(n)}; \mathbf{z}) k(\mathbf{z}^{(m)}; \mathbf{z}') \\ &= \frac{F}{P^2} \sum_n \left(\frac{\lambda^{(n)}}{\mathbb{E}_t[\lambda_t]} \right)^2 \text{Cov} \left[\tilde{\mathbf{s}}(\mathbf{z}^{(n)}; \mathbf{x}_0^{(n)}) \right] k(\mathbf{z}^{(n)}; \mathbf{z}) k(\mathbf{z}^{(n)}; \mathbf{z}') \\ &= \frac{F}{P} \int \frac{\lambda_{t''}^2}{\mathbb{E}_t[\lambda_t]^2} \mathbf{C}(\mathbf{z}'') k(\mathbf{z}''; \mathbf{z}) k(\mathbf{z}''; \mathbf{z}') p(\mathbf{z}'') d\mathbf{z}'' ,\end{aligned}$$

when P is large, where we exploited the independence of the samples in the first step, and the central limit theorem in the second. As elsewhere, we have used $\mathbf{C}(\mathbf{z}) := \text{Cov}[\tilde{\mathbf{s}}(\mathbf{z}; \mathbf{x}_0)]$ as shorthand.

Finally, the V-kernel is

$$\mathbf{V}(\mathbf{z}; \mathbf{z}') = \lim_{P \rightarrow \infty} \frac{F}{P} \mathbf{D}_t \mathbb{E}_{\mathbf{z}''} \left\{ \frac{\lambda_{t''}^2}{\mathbb{E}[\lambda_t]^2} k(\mathbf{z}; \mathbf{z}'') \mathbf{C}(\mathbf{z}'') k(\mathbf{z}''; \mathbf{z}') \right\} \mathbf{D}_{t'} \quad (101)$$

provided that the limit exists and is finite. Since $F \propto N$, this can happen if $N \propto P$.

Also note that the form of this V-kernel is identical to that of the V-kernel for linear models (c.f. Prop. 5.1), with the only difference being the feature kernel that appears.

¹We have run out of letters, and unfortunately will use k to denote this quantity. Note that it is different from both \mathbf{K} and \mathbf{k} .

The infinite training time limit is of particular interest, since in this limit we expect the model to interpolate all of its (noisy) samples. In this limit, we have

$$\begin{aligned}
\text{Cov}[\hat{s}_i(\mathbf{z}), \hat{s}_j(\mathbf{z}')] &= \frac{1}{P} \int \frac{\lambda_{t''}^2}{\mathbb{E}_t[\lambda_t]^2} \phi(\mathbf{z})^T \phi(\mathbf{z}'') \phi(\mathbf{z}'')^T \phi(\mathbf{z}') C_{ij}(\mathbf{z}'') p(\mathbf{z}'') d\mathbf{z}'' \\
&= \frac{1}{P} \int \frac{\lambda_{t''}}{\mathbb{E}_t[\lambda_t]} \phi(\mathbf{z})^T \phi(\mathbf{z}'') \left[\frac{\lambda_{t''}}{\mathbb{E}_t[\lambda_t]} \phi(\mathbf{z}'')^T \phi(\mathbf{z}') p(\mathbf{z}'') \right] C_{ij}(\mathbf{z}'') d\mathbf{z}'' \\
&= \frac{1}{P} \int \frac{\lambda_{t''}}{\mathbb{E}_t[\lambda_t]} \phi(\mathbf{z})^T \phi(\mathbf{z}'') \delta(\mathbf{z}'' - \mathbf{z}') C_{ij}(\mathbf{z}'') d\mathbf{z}'' \\
&= \frac{1}{P} \frac{\lambda_{t'}}{\mathbb{E}_t[\lambda_t]} \phi(\mathbf{z})^T \phi(\mathbf{z}') C_{ij}(\mathbf{z}')
\end{aligned} \tag{102}$$

where we have exploited the completeness relation. Now we encounter a subtle technical point. Since $F \neq P$ in general, in the $F, P \rightarrow \infty$ limit the quantity

$$d(\mathbf{z}, \mathbf{z}') := \frac{1}{P} \frac{\lambda_{t'}}{\mathbb{E}_t[\lambda_t]} \Phi(\mathbf{z})^T \Phi(\mathbf{z}') \tag{103}$$

is not quite equal to the Dirac delta function, but is instead proportional to it. We need to work out the constant of proportionality. To do this, observe that

$$\sum_n d(\mathbf{z}^{(n)}, \mathbf{z}^{(n)}) = \frac{1}{P} \sum_n \frac{\lambda^{(n)}}{\mathbb{E}_t[\lambda_t]} \Phi(\mathbf{z}^{(n)})^T \Phi(\mathbf{z}^{(n)}) \rightarrow \sum_{k=1}^F \int \frac{\lambda_t}{\mathbb{E}_t[\lambda_t]} \phi_k(\mathbf{z}) \phi_k(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} = F.$$

On the other hand, for the Dirac delta function, we would have

$$\sum_n \delta(0) = \frac{P}{\Delta \mathbf{z}}, \tag{104}$$

where $\Delta \mathbf{z}$ is some small bin size. This implies

$$d(\mathbf{z}, \mathbf{z}') = \frac{F \Delta \mathbf{z}}{P} \delta(\mathbf{z} - \mathbf{z}'). \tag{105}$$

If we define $\kappa := (F \Delta \mathbf{z})/P$, and assume κ remains constant as both parameters approach infinity, we finally obtain

$$\mathbf{V}(\mathbf{z}; \mathbf{z}') = \kappa \mathbf{D}_t \mathbf{C}(\mathbf{z}) \mathbf{D}_t \delta(\mathbf{z} - \mathbf{z}'). \tag{106}$$

H RESULTS FOR NOISE PREDICTION FORMULATION

In the main text, we present the problem of training a diffusion model in terms of learning a parameterized score function estimator (see Eq. 4). In practice, one often does not try to directly estimate the score, but the quantity $\mathcal{E}(\mathbf{x}, t) := \sigma_t \mathbf{s}(\mathbf{x}, t)$ (Rombach et al., 2022), since it can be somewhat better-behaved. In this appendix, we note that it is straightforward to port our results to this slightly different setting.

For simplicity, assume that the forward process involves isotropic noise, so that $\mathbf{D}_t = \frac{g_t^2}{2} \mathbf{I}_D$ and $\mathbf{S}_t = \sigma_t^2 \mathbf{I}_D$. The DSM objective (Eq. 4) becomes

$$\begin{aligned} J_1(\boldsymbol{\theta}) &= \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}} \left\{ \frac{\lambda_t}{2\sigma_t^2} \|\sigma_t \hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \sigma_t \bar{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)\|_2^2 \right\} \\ &= \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}} \left\{ \frac{\lambda_t}{2\sigma_t^2} \|\hat{\mathcal{E}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathcal{E}\|_2^2 \right\}. \end{aligned} \quad (107)$$

Since $\mathcal{E} := (\alpha_t \mathbf{x}_0 - \mathbf{x})/\sigma_t$ by definition, \mathcal{E} is normally distributed with mean zero and variance one. The function $\hat{\mathcal{E}}$ can be viewed as taking a noisy sample \mathbf{x} as input and outputting the (standardized) noise that was added to the original sample \mathbf{x}_0 .

Importantly, the objective has the same form as before (i.e., a mean-squared error objective comparing an estimator to a target), but the time-weighting function is now $\tilde{\lambda}_t := \lambda_t/\sigma_t^2$.

In this formulation, the PF-ODE reads

$$\begin{aligned} \dot{\mathbf{x}}_t &= -\beta_t \mathbf{x}_t - \frac{g_t^2}{2\sigma_t} \mathcal{E}(\mathbf{x}, t) \\ &= -\beta_t \mathbf{x}_t - \tilde{D}_t \mathcal{E}(\mathbf{x}, t) \end{aligned} \quad (108)$$

where we define

$$\tilde{D}_t := \frac{g_t^2}{2\sigma_t}. \quad (109)$$

Hence, the PF-ODE also has the same form as before. We conclude that our results apply once one makes these identifications, and also makes the slight change

$$\mathbf{C}(\mathbf{x}, t) := \sigma_t^2 \text{Cov}_{\mathbf{x}_0|\mathbf{x}, t}(\bar{\mathbf{s}}) \quad (110)$$

since the learning target is now a scaled version $\sigma_t \bar{\mathbf{s}}$ of the proxy score.

I RESULTS FOR DENOISER FORMULATION

In the main text, we present the problem of training a diffusion model in terms of learning a parameterized score function estimator (see Eq. 4). An alternative approach formulates the problem in terms of learning a ‘denoiser’ function, which takes a noise-corrupted sample \mathbf{x} as input and outputs a noise-free sample \mathbf{x}_0 . These formulations are mathematically equivalent, but their implementations slightly differ in practice; in this appendix, we note that it is straightforward to port our results to this slightly different setting.

Define the ‘optimal’ denoiser² \mathbf{D} in terms of the true score \mathbf{s} via

$$\mathbf{D}(\mathbf{x}, t) := \frac{1}{\alpha_t} [\mathbf{S}_t \mathbf{s}(\mathbf{x}, t) + \mathbf{x}] \quad \mathbf{s}(\mathbf{x}, t) = \mathbf{S}_t^{-1} [\alpha_t \mathbf{D}(\mathbf{x}, t) - \mathbf{x}] . \quad (111)$$

Parameterized estimators of the denoiser and score, which we will denote by $\hat{\mathbf{D}}$ and $\hat{\mathbf{s}}$, are related in the same way. This means that the DSM objective (Eq. 4) becomes

$$\begin{aligned} J_1(\boldsymbol{\theta}) &= \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}} \left\{ \frac{\lambda_t}{2} \|\hat{\mathbf{s}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \tilde{\mathbf{s}}(\mathbf{x}, t; \mathbf{x}_0)\|_2^2 \right\} \\ &= \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}} \left\{ \frac{\lambda_t}{2} \|\mathbf{S}_t^{-1} [\alpha_t \hat{\mathbf{D}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{x}] - \mathbf{S}_t^{-1} [\alpha_t \mathbf{x}_0 - \mathbf{x}]\|_2^2 \right\} \\ &= \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}} \left\{ \frac{\lambda_t \alpha_t^2}{2} [\hat{\mathbf{D}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{x}_0]^T \mathbf{S}_t^{-2} [\hat{\mathbf{D}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{x}_0] \right\} . \end{aligned} \quad (112)$$

For simplicity, we will assume the diffusion tensor of the forward process is isotropic ($\mathbf{D}_t = \frac{g_t^2}{2} \mathbf{I}_D$), which implies $\mathbf{S}_t = \sigma_t^2 \mathbf{I}_D$. Then

$$J_1(\boldsymbol{\theta}) = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}} \left\{ \frac{\lambda_t \alpha_t^2}{2 \sigma_t^4} \|\hat{\mathbf{D}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{x}_0\|_2^2 \right\} = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}} \left\{ \frac{\tilde{\lambda}_t}{2} \|\hat{\mathbf{D}}_{\boldsymbol{\theta}}(\mathbf{x}, t) - \mathbf{x}_0\|_2^2 \right\} \quad (113)$$

where $\tilde{\lambda}_t := \lambda_t \alpha_t^2 / \sigma_t^4$. Hence, the objective has the same form as before (a mean-squared error objective comparing an estimator to a target, with a particular time-weighting function).

In this formulation, the PF-ODE reads

$$\begin{aligned} \dot{\mathbf{x}}_t &= -\beta_t \mathbf{x}_t - \frac{g_t^2}{2\sigma_t^2} [\alpha_t \mathbf{D}(\mathbf{x}_t, t) - \mathbf{x}] \\ &= -\left(\beta_t - \frac{g_t^2}{2\sigma_t^2} \right) \mathbf{x}_t - \frac{g_t^2}{2\sigma_t^2} \alpha_t \mathbf{D}(\mathbf{x}_t, t) \\ &= -\tilde{\beta}_t \mathbf{x}_t - \tilde{D}_t \mathbf{D}(\mathbf{x}_t, t) \end{aligned} \quad (114)$$

where we define

$$\tilde{\beta}_t := \beta_t - \frac{g_t^2}{2\sigma_t^2} \quad \tilde{D}_t := \frac{g_t^2}{2\sigma_t^2} \alpha_t . \quad (115)$$

Hence, the PF-ODE also has the same form as before. We conclude that our results apply once one makes these identifications, and also makes the slight change

$$\mathbf{C}(\mathbf{x}, t) := \text{Cov}_{\mathbf{x}_0 | \mathbf{x}, t}(\mathbf{x}_0) \quad (116)$$

since the learning target is now \mathbf{x}_0 rather than the proxy score.

²There is an unfortunate collision of notation here, with \mathbf{D} being used to denote both the diffusion tensor and the denoiser, but it should be fairly clear from context which is which.

J BENIGN PROPERTIES OF GENERALIZATION THROUGH VARIANCE

A priori, one may worry that generalization through variance can happen in an ‘uncontrolled’ fashion, and hence produce generalizations of the data distribution that are somehow ‘bad’ or ‘unreasonable’. In this appendix, we collect properties of generalization through variance that help ensure that this does not happen.

Most of these properties relate to the proxy score covariance, which when the data distribution consists of M examples has the form (see Appendix B)

$$\mathbf{C}(\mathbf{x}, t) := \text{Cov}_{\mathbf{x}_0|\mathbf{x},t}(\tilde{\mathbf{s}}) = \alpha_t^2 \mathbf{S}_t^{-1} \text{Cov}_{\mathcal{M}}(\boldsymbol{\mu}) \mathbf{S}_t^{-1} \quad (117)$$

where

$$p(\mathbf{x}_0 = \boldsymbol{\mu}_m | \mathbf{x}, t) = \frac{\mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_m, \mathbf{S}_t)}{\sum_{m'} \mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_{m'}, \mathbf{S}_t)}. \quad (118)$$

For the sake of this appendix, we will focus mainly on the case where the forward process is isotropic, which means $\mathbf{G}_t = g_t \mathbf{I}_D$, $\mathbf{D}_t = \frac{g_t^2}{2} \mathbf{I}_D$, and $\mathbf{S}_t = \sigma_t^2 \mathbf{I}_D$. In this case,

$$\begin{aligned} \mathbf{C}(\mathbf{x}, t) &= \frac{\alpha_t^2}{\sigma_t^4} \text{Cov}_{\mathbf{x}_0|\mathbf{x},t}(\mathbf{x}_0) \\ p(\mathbf{x}_0 = \boldsymbol{\mu}_m | \mathbf{x}, t) &= \frac{\mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_m, \sigma_t^2 \mathbf{I}_D)}{\sum_{m'} \mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}_{m'}, \sigma_t^2 \mathbf{I}_D)} = \text{softmax} \left(-\frac{\|\mathbf{x} - \alpha_t \boldsymbol{\mu}_m\|_2^2}{2\sigma_t^2} \right). \end{aligned} \quad (119)$$

J.1 SINGLE POINTS ARE NOT GENERALIZED

Suppose that the data distribution consists of a single point at $\mathbf{x}_0 = \boldsymbol{\mu}$, so that $M = 1$. How do diffusion models generalize a single point? Since

$$p(\mathbf{x}_0 = \boldsymbol{\mu} | \mathbf{x}, t) = \frac{\mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}, \sigma_t^2 \mathbf{I}_D)}{\mathcal{N}(\mathbf{x}; \alpha_t \boldsymbol{\mu}, \sigma_t^2 \mathbf{I}_D)} = 1, \quad (120)$$

the covariance of \mathbf{x}_0 is zero. This makes sense, since given a pair (\mathbf{x}, t) , there is no uncertainty about the training example \mathbf{x}_0 that generated it.

This implies that the covariance of \mathbf{x}_0 given \mathbf{x} and t is zero, and hence (if the forward process is isotropic) that the proxy score covariance is zero. For the naive score estimator and NTK-regime neural network, this implies that the V-kernel is always zero, and hence that generalization through variance does not occur. This is reasonable, among other reasons because there is no reason to introduce anisotropy if it is not present in the data.

J.2 DIMENSIONALITY OF DATA DISTRIBUTION IS PRESERVED

Suppose that the data distribution lies entirely within a subspace of dimension $r < D$, so that the components of each $\boldsymbol{\mu}_m$ along the non-subspace dimensions are equal (to μ'_k along dimension k , say). It seems intuitively reasonable not to substantially modify probability mass outside of this subspace. We clearly have

$$\mathbb{E}_{\mathbf{x}_0|\mathbf{x},t}[x_{0k}] = \mu'_k \quad (121)$$

for all non-subspace directions k . This has the following consequence. Let ℓ denote some other (possibly within the subspace, possibly not) direction of state space. The k - ℓ covariance is

$$\text{Cov}_{\mathbf{x}_0|\mathbf{x},t}[(x_{0k} - \mu'_k)(x_{0\ell} - \mathbb{E}_{\mathbf{x}_0|\mathbf{x},t}[x_{0\ell}])] = 0 \quad (122)$$

since x_{0k} is always equal to μ'_k . Hence, C_{ij} (and by extension, V_{ij}) is only nonzero if i and j are directions along which the data distribution varies, which means generalization through variance only happens along the ‘data manifold’.

J.3 VARIANCE IS NOT ADDED FAR FROM DATA DISTRIBUTION EXAMPLES

Since $p(\mathbf{x}_0|\mathbf{x}, t)$ has the form of a softmax function, taking \mathbf{x} extremely large is analogous to taking the temperature parameter of a typical softmax to be extremely small. For \mathbf{x} far from the support

of the data distribution, to good approximation $p(\mathbf{x} = \boldsymbol{\mu}_m | \mathbf{x}, t) = \delta_{m, m_C(\mathbf{x}, t)}$, where $m_C(\mathbf{x}, t)$ denotes the index of the data point closest (in terms of Euclidean distance) to \mathbf{x} . Since $p(\mathbf{x}_0 | \mathbf{x}, t)$ collapses to a Kronecker delta function, its covariance goes to zero, and hence the V-kernel also goes to zero.

J.4 FOLLOWING AVERAGE SCORE FIELD IS MOST LIKELY

A straightforward consequence of the effective reverse process (Eq. 8) is that

$$\frac{d}{dt} \langle \mathbf{x} \rangle = \langle \mathbf{f}(\mathbf{x}, t) \rangle \quad (123)$$

where the angular brackets here denote averages with respect to $[q(\mathbf{x}, t)]$. Hence, although the effective reverse process involves noise, that noise has zero mean, so on average the system follows PF-ODE dynamics that use the ensemble-averaged score estimator. If the estimator is unbiased, this is just the usual PF-ODE. In other words, even in our setting, *on average* PF-ODE dynamics will reproduce training examples.

There is a slight technical subtlety here, which is the fact that $\langle \mathbf{f}(\mathbf{x}, t) \rangle \neq \mathbf{f}(\langle \mathbf{x} \rangle, t)$ in general, but this distinction becomes unimportant for very small noise scales, which most strongly influence whether memorization occurs.

J.5 TRAINING DATA ARE MORE LIKELY TO BE SAMPLED WHEN NOISE IS SMALL

Similar to the previous point, if the prefactor κ that controls the size of the V-kernel is not too large, trajectories that follow PF-ODE dynamics are more *likely* than trajectories that do not. Note that there is a distinction between the *most likely* trajectories of a stochastic process, and the *average* trajectory associated with that process, although here one expects them to be fairly similar.

This idea can be formalized in the context of a semiclassical analysis (see Appendix ??), which shows that the probability of a given path goes like $\exp(-S/\kappa)$, where

$$\mathcal{S}[\mathbf{x}_t] := \int_{\epsilon}^T \int_{\epsilon}^T \frac{1}{2} [\dot{\mathbf{x}}_t - \mathbf{f}(\mathbf{x}_t, t)]^T \mathbf{Q}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') [\dot{\mathbf{x}}_{t'} - \mathbf{f}(\mathbf{x}_{t'}, t')] dt dt' \quad (124)$$

for some matrix \mathbf{Q} that functions as the inverse of \mathbf{V} . The point we would like to make here corresponds to the observation that this action attains its smallest value when $\dot{\mathbf{x}}_t = \mathbf{f}(\mathbf{x}_t, t)$ for all times t . That is, trajectories which follow PF-ODE dynamics are more likely than those that deviate from it, with deviations being penalized more harshly as κ is made smaller.

K MEMORIZATION AND THE V-KERNEL IN THE SMALL NOISE LIMIT

It's clear how the V-kernel affects a given sampling time step, but how does it affect the overall learned distribution? This question is difficult to answer analytically, since in general the effective reverse process is not exactly solvable. Even without the difficulties related to state-dependent noise and nontrivial temporal autocorrelations that the V-kernel introduces, PF-ODE dynamics are only exactly solvable in special circumstances, e.g., when the data distribution is Gaussian (Wang & Vastola, 2023; 2024).

It is slightly easier to relate the V-kernel to the (average) learned distribution $[q(\mathbf{x}_0)]$ when the overall magnitude of the V-kernel is small, since in this limit one can invoke a semiclassical approximation (Kleinert, 2006), which is the path integral analogue of a saddlepoint approximation. This limit plausibly applies if the prefactor κ (which is proportional to the ratio F/P for the latter two kinds of models we consider) is small, or equivalently if the model is somewhat underparameterized.

A proper treatment of this topic involves the careful manipulation of notoriously tricky mathematical objects like functional determinants; to sidestep these issues, and be somewhat more brief, we will proceed in a somewhat heuristic fashion here. If one of the following steps appears unclear, we advise the reader to examine it in discrete rather than continuous time, where the associated issues are less severe.

K.1 SETTING UP THE SEMICLASSICAL APPROXIMATION

First, define the prefactor-divided V-kernel \tilde{V} as

$$\tilde{V}(\mathbf{z}; \mathbf{z}') := \mathbf{V}(\mathbf{z}; \mathbf{z}') / \kappa, \quad (125)$$

and define the ‘inverse’ \mathbf{Q} of the V-kernel as the vector-valued matrix satisfying

$$\int_{\epsilon}^T \mathbf{Q}(\mathbf{x}_t, t; \mathbf{x}_{t''}, t'') \tilde{V}(\mathbf{x}_{t''}, t''; \mathbf{x}_{t'}, t') dt'' = \int_{\epsilon}^T \tilde{V}(\mathbf{x}_t, t; \mathbf{x}_{t''}, t'') \mathbf{Q}(\mathbf{x}_{t''}, t''; \mathbf{x}_{t'}, t') dt'' = \mathbf{I}_D \delta(t-t').$$

Our path integral expression for $[q(\mathbf{x}_0)]$ has the form (see Appendix D)

$$[q(\mathbf{x}_0)] = \int \mathcal{D}[\mathbf{x}_t] \mathcal{D}[\mathbf{p}_t] \exp\{-\mathcal{S}[\mathbf{x}_t, \mathbf{p}_t]\} \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{S}_T), \quad (126)$$

where the ‘action’ \mathcal{S} is

$$\mathcal{S}[\mathbf{x}_t, \mathbf{p}_t] := \int_{\epsilon}^T -i \mathbf{p}_t^T [\dot{\mathbf{x}}_t - \mathbf{f}(\mathbf{x}_t, t)] dt + \frac{\kappa}{2} \int_{\epsilon}^T \int_{\epsilon}^T \mathbf{p}_t^T \tilde{V}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') \mathbf{p}_{t'} dt dt', \quad (127)$$

and where the functional integral is over all paths $\mathbf{x}(t)$ which have $\mathbf{x}(0) = \mathbf{x}_0$ and $\mathbf{x}(T) = \mathbf{x}_T$, and all possible $\mathbf{p}(t)$ paths. The explicit form of \mathbf{f} is

$$\mathbf{f}(\mathbf{x}, t) := -\beta_t \mathbf{x} - \mathbf{D}_t \mathbf{s}_{avg}(\mathbf{x}, t), \quad (128)$$

although for this analysis it is not relevant.

Since the action is quadratic in the ‘momenta’ variables \mathbf{p}_t , we can perform the associated Gaussian integrals exactly to obtain a reduced path integral with

$$[q(\mathbf{x}_0)] = \int \mathcal{D}[\mathbf{x}_t] \mathcal{D}[\mathbf{p}_t] \exp\{-\mathcal{S}[\mathbf{x}_t] / \kappa\} F(\{\mathbf{x}_t\}) \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{S}_T) \quad (129)$$

$$\mathcal{S}[\mathbf{x}_t] = \int_{\epsilon}^T \int_{\epsilon}^T [\dot{\mathbf{x}}_t - \mathbf{f}(\mathbf{x}_t, t)]^T \mathbf{Q}(\mathbf{x}_t, t; \mathbf{x}_{t'}, t') [\dot{\mathbf{x}}_{t'} - \mathbf{f}(\mathbf{x}_{t'}, t')] dt dt'$$

where the factor F involves a functional determinant

$$F(\{\mathbf{x}_t\}) := \frac{1}{\sqrt{\det \mathbf{V}}} \quad (130)$$

due to the Hessian of the action. Note that \mathbf{Q} is positive semidefinite since \mathbf{V} is, and hence the minimum possible value of \mathcal{S} is zero.

K.2 SEMICLASSICAL APPROXIMATION OF THE LEARNED DISTRIBUTION

One can now invoke the semiclassical approximation of this path integral assuming κ is sufficiently small. Relevant to this approximation is the ‘classical’ action $\mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T)$, which quantifies the likelihood of \mathbf{x}_t following the most likely (i.e., ‘classical’) path from \mathbf{x}_T to \mathbf{x}_0 :

$$\mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T) := \min_{\mathbf{x}(t): \mathbf{x}(0)=\mathbf{x}_0, \mathbf{x}(T)=\mathbf{x}_T} \mathcal{S}[\mathbf{x}_t]. \quad (131)$$

Also relevant is the Hessian of this quantity evaluated at the least action path, which reads

$$\mathcal{H} = \frac{1}{\kappa} \frac{\delta^2 \mathcal{S}}{\delta \mathbf{x}_t \delta \mathbf{x}_{t'}} = \frac{1}{\kappa} \left(\frac{d}{dt} - \mathbf{J}_f \right)^T \mathbf{Q} \left(\frac{d}{dt} - \mathbf{J}_f \right) + \text{additional terms} \quad (132)$$

where \mathbf{J}_f is the Jacobian of \mathbf{f} . The additional terms will vanish after we make one more approximation, so we ignore them. The (functional) determinant of this Hessian is

$$\det \mathcal{H}(\mathbf{x}_0, \mathbf{x}_T) = \det \left[\frac{1}{\kappa} \left(\frac{d}{dt} - \mathbf{J}_f \right)^T \mathbf{Q} \left(\frac{d}{dt} - \mathbf{J}_f \right) + \text{additional terms} \right]. \quad (133)$$

Although our notation here is deliberately somewhat vague to sidestep various technical details, the above quantity depends on the entire path $\{\mathbf{x}_t\}$, and the values of the Jacobian and \mathbf{Q} along it.

The semiclassical approximation says that

$$[q(\mathbf{x}_0)] \approx \int \frac{d\mathbf{x}_T}{\sqrt{(2\pi)^D}} \exp \{ -\mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T)/\kappa \} \frac{\mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{S}_T)}{\sqrt{\det \mathcal{H}(\mathbf{x}_0, \mathbf{x}_T) \det \mathbf{V}(\mathbf{x}_0, \mathbf{x}_T)}}. \quad (134)$$

We can invoke Laplace’s method in order to approximately evaluate the \mathbf{x}_T integral, and hence obtain an expression for the (average) learned distribution that does not involve any integrals. The classical action \mathcal{S} can be expanded with respect to $\mathbf{x}_T^*(\mathbf{x}_0)$, the *most likely* (in terms of minimizing \mathcal{S}_{cl}) noise seed \mathbf{x}_T given the endpoint \mathbf{x}_0 :

$$\mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T) \approx \mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0)) + \frac{1}{2} [\mathbf{x}_T - \mathbf{x}_T^*]^T \frac{\partial^2 \mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0))}{\partial \mathbf{x}_T \partial \mathbf{x}_T} [\mathbf{x}_T - \mathbf{x}_T^*]. \quad (135)$$

But the classical action takes its minimum possible value—i.e., zero—when \mathbf{x}_T is chosen to be the unique noise seed that corresponds to deterministic PF-ODE dynamics (i.e., $\dot{\mathbf{x}}_t = \mathbf{f}(\mathbf{x}_t, t)$ at all times t), so we just have

$$\mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T) \approx \frac{1}{2} [\mathbf{x}_T - \mathbf{x}_T^*]^T \frac{\partial^2 \mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0))}{\partial \mathbf{x}_T \partial \mathbf{x}_T} [\mathbf{x}_T - \mathbf{x}_T^*]. \quad (136)$$

By Laplace’s method, which is also usable due to the smallness of κ , we obtain

$$[q(\mathbf{x}_0)] \approx \frac{1}{\sqrt{\det \left(\frac{1}{\kappa} \frac{\partial^2 \mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0))}{\partial \mathbf{x}_T \partial \mathbf{x}_T} \right)}} \frac{\mathcal{N}(\mathbf{x}_T^*(\mathbf{x}_0); \mathbf{0}, \mathbf{S}_T)}{\sqrt{\det \mathcal{H}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0)) \det \mathbf{V}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0))}}. \quad (137)$$

If we only consider ‘classical’ paths with $\dot{\mathbf{x}}_t = \mathbf{f}(\mathbf{x}_t, t)$, as in this approximation, the additional terms in Eq. 132 vanish. This means

$$\begin{aligned} \det \mathcal{H}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0)) &= \det \left[\frac{1}{\kappa} \left(\frac{d}{dt} - \mathbf{J}_f \right)^T \mathbf{Q} \left(\frac{d}{dt} - \mathbf{J}_f \right) \right] \\ &= \det \left(\frac{\mathbf{Q}}{\kappa} \right) \det \left(\frac{d}{dt} - \mathbf{J}_f \right)^2 \end{aligned} \quad (138)$$

where this manipulation can be more formally justified if one works in discrete time. But since \mathbf{Q}/κ is the inverse of \mathbf{V} , and since

$$\det \left(\frac{d}{dt} - \mathbf{J}_f \right) \approx \det \left(\prod_{t=1}^T [\mathbf{I} - \mathbf{J}_t \Delta t] \right) \approx \det \exp \left\{ - \int_{\epsilon}^T \mathbf{J}_f(t) dt \right\}, \quad (139)$$

where we have used a discrete time argument to compute the determinant, we have

$$\begin{aligned} \det \mathcal{H}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0)) \det \mathbf{V}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0)) &= \det \mathbf{V} \det \mathbf{V}^{-1} \det \exp \left\{ -2 \int_{\epsilon}^T \mathbf{J}_{\mathbf{f}}(t) dt \right\} \\ &= \det \exp \left\{ -2 \int_{\epsilon}^T \mathbf{J}_{\mathbf{f}}(t) dt \right\}. \end{aligned} \quad (140)$$

Note that this determinant can be simplified somewhat:

$$\log \det \exp \left\{ -2 \int_{\epsilon}^T \mathbf{J}_{\mathbf{f}}(t) dt \right\} = -2 \int_{\epsilon}^T \text{tr}[\mathbf{J}_{\mathbf{f}}(t)] dt = -2 \int_{\epsilon}^T \nabla_{\mathbf{x}_t} \cdot \mathbf{f}(\mathbf{x}_t, t) dt \quad (141)$$

where the \mathbf{x}_t that appears in the expression above follows deterministic PF-ODE dynamics. Finally,

$$[q(\mathbf{x}_0)] \approx \frac{1}{\sqrt{\det \left(\frac{1}{\kappa} \frac{\partial^2 \mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0))}{\partial \mathbf{x}_T \partial \mathbf{x}_T} \right)}} \exp \left\{ \log \mathcal{N}(\mathbf{x}_T^*(\mathbf{x}_0); \mathbf{0}, \mathbf{S}_T) + \int_{\epsilon}^T \nabla_{\mathbf{x}_t} \cdot \mathbf{f}(\mathbf{x}_t, t) dt \right\}. \quad (142)$$

At this point, we can make a crucial observation: the argument of the exponential is *precisely* the instantaneous change of variables formula that can be used to compute the log-likelihood $p(\mathbf{x}_0|\epsilon)$ (Song et al., 2021; Chen et al., 2018). This immediately implies

$$[q(\mathbf{x}_0)] \approx p(\mathbf{x}_0|\epsilon) \frac{1}{\sqrt{\det \left(\frac{1}{\kappa} \frac{\partial^2 \mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0))}{\partial \mathbf{x}_T \partial \mathbf{x}_T} \right)}}. \quad (143)$$

We conclude that, at least in the small κ regime, the learned distribution is equal to the memorized (ϵ -noise-corrupted) data distribution, times the determinant of a Hessian that quantifies the likelihood of deviating from PF-ODE dynamics. This Hessian depends on the V-kernel, since the classical action depends on its inverse \mathbf{Q} , so it is precisely here that the V-kernel can influence generalization.

The required Hessian appears difficult to compute in general. Incidentally, since the relevant action (Eq. 129) is generically not local in time, it is also hard to derive a Hamilton-Jacobi-type differential equation satisfied by this Hessian.

K.3 QUANTIFYING MEMORIZATION IN THE SEMICLASSICAL REGIME

Suppose we quantify memorization by computing the Kullback-Leibler (KL) divergence between the data distribution p_{data} and the (average) learned distribution $[q]$:

$$E_{mem} := D_{KL}(p_{data} \parallel [q]) = \int p_{data}(\mathbf{x}_0) \log \frac{p_{data}(\mathbf{x}_0)}{[q(\mathbf{x}_0)]} d\mathbf{x}_0. \quad (144)$$

By our semiclassical approximation result (Eq. 143), this is just

$$E_{mem} = D_{KL}(p_{data} \parallel p_{\epsilon}) + \frac{1}{2} \int p_{data}(\mathbf{x}_0) \log \det \left(\frac{1}{\kappa} \frac{\partial^2 \mathcal{S}_{cl}(\mathbf{x}_0, \mathbf{x}_T^*(\mathbf{x}_0))}{\partial \mathbf{x}_T \partial \mathbf{x}_T} \right) d\mathbf{x}_0, \quad (145)$$

where $p_{\epsilon} := p(\mathbf{x}_0|\epsilon)$. Hence, the curvature of the classical action near $\mathbf{x}_T^*(\mathbf{x}_0)$ strongly controls the extent to which the data distribution is memorized, especially when ϵ is taken to be small, in which case the first term is negligible.