# A Complete Recipe for Diffusion Generative Models

**Kushagra Pandey**
Department of Computer Science
University of California, Irvine
pandeyk1@uci.edu

**Stephan Mandt**
Department of Computer Science
University of California, Irvine
mandt@uci.edu

## Abstract

Score-based Generative Models (SGMs) have demonstrated exceptional synthesis outcomes across various tasks. However, the current design landscape of the forward diffusion process remains largely untapped and often relies on physical heuristics or simplifying assumptions. Utilizing insights from the development of scalable Bayesian posterior samplers, we present a complete recipe for formulating forward processes in SGMs, ensuring convergence to the desired target distribution. Our approach reveals that several existing SGMs can be seen as specific manifestations of our framework. Building upon this method, we introduce Phase Space Langevin Diffusion (PSLD), which relies on score-based modeling within an augmented space enriched by auxiliary variables akin to physical phase space. Empirical results exhibit the superior sample quality and improved speed-quality trade-off of PSLD compared to various competing approaches on established image synthesis benchmarks. Remarkably, PSLD achieves sample quality akin to state-of-the-art SGMs (FID: **2.10** for unconditional CIFAR-10 generation). Lastly, we demonstrate the applicability of PSLD in conditional synthesis using pre-trained score networks, offering an appealing alternative as an SGM backbone for future advancements. Code and model checkpoints can be accessed at `https://github.com/mandt-lab/PSLD`.

## 1 Introduction

Score-based Generative Models [1–4] are a class of explicit-likelihood based generative models that have recently demonstrated impressive performance on various synthesis benchmarks, such as image generation [5–9], video synthesis [10, 11] and 3D shape generation [12, 13]. SGMs employ a forward stochastic process to add noise to data incrementally, transforming the data-generating distribution to a tractable prior distribution that enables sampling. Subsequently, a learnable reverse process transforms the prior distribution back to the data distribution using a parametric estimator of the gradient field of the log probability density of the data (a.k.a score).

However, a principled framework for extending the current design space of diffusion processes is still missing. Although some studies have proposed augmenting the forward diffusion process with auxiliary variables [14] to improve sample quality, their design is primarily motivated by physical intuition and non-obvious how to generalize. Therefore, a principled framework is required to explore the space of possible diffusion processes better.

In this work, we propose a *complete* recipe for the design of diffusion processes, motivated by the design of stochastic gradient MCMC samplers [15–17]. Our recipe leads to a flexible parameterization of the forward diffusion process without requiring physical intuition. Moreover, under the proposed parameterization, the forward process is guaranteed to converge to a prior distribution of interest. We show that several existing SGMs can be cast under our diffusion process parameterization. Furthermore, using our proposed recipe, we introduce PSLD, a novel SGM which performs diffusion in the joint space of data and auxiliary variables. We demonstrate that PSLD generalizes Critically

Figure 1: Unconditional PSLD generated samples. AFHQv2 128 x 128 (Top), CelebA-64 (Bottom Left, FID=2.01) and CIFAR-10 (Bottom Right, FID=2.10)

Damped Langevin Diffusion (CLD) [14] and outperforms existing baselines on several empirical settings on standard image synthesis benchmarks such as CIFAR-10 [18] and CelebA-64 [19]. More specifically, we make the following theoretical and empirical contributions:

1. **A Complete Recipe for SGM Design:** We propose a specific parameterization of the forward process, guaranteed to converge asymptotically to a desired stationary "prior" distribution. The proposed recipe is *complete* in the sense that it subsumes all possible Markovian stochastic processes which converge to this distribution. We show that several existing SGMs [4, 14] can be cast as specific instantiations of our recipe.

2. **Phase Space Langevin Diffusion(PSLD):** To exemplify the proposed diffusion parameterization concretely, we propose PSLD: a novel SGM which performs diffusion in the phase space by adding noise in **both** data and the momentum space.

3. **Superior Sample Quality and Speed-Quality Tradeoffs:** Using ablation experiments on standard image synthesis benchmarks like CIFAR-10 and CelebA-64, we demonstrate the benefits of adding stochastic noise in both the data and the momentum space on overall sample quality and the speed quality trade-offs associated with PSLD. Furthermore, using similar score network architectures, our proposed method outperforms existing diffusion baselines on both criteria across different sampler settings.

4. **State-of-the-Art Sample Quality:** We show that PSLD outperforms competing baselines and achieves competitive perceptual sample quality to other state-of-the-art methods. Our model achieves an FID [20] score of **2.10**, an IS score [21] of **9.93** on unconditional CIFAR-10 and an FID score of **2.01** on CelebA-64.

5. **Conditional synthesis:** We show that pre-trained unconditional PSLD models can be used for conditional synthesis tasks like class-conditional generation and image inpainting.

Overall, given the superior performance of PSLD on several tasks, we present an attractive alternative to existing SGM backbones for further development. We organize the rest of our work as follows: Section 2 presents some background on SGMs and our proposed *recipe* for SGM design. Section 3 presents the construction of our novel PSLD model. Section 4 presents our empirical findings. Lastly, Section 5 compares our proposed contributions to several existing works while we present some directions for future work in Section 6.

## 2 A Complete Recipe for SGM Design

### 2.1 Background

Consider the following forward process SDE for converting data $\mathbf{x}_t \in \mathbb{R}^d$ to noise,

$$d\mathbf{x}_t = \boldsymbol{f}(\mathbf{x}_t, t)\, dt + \boldsymbol{G}(t)\, d\mathbf{w}_t, \quad t \in [0, T],$$

with continuous time variable $t \in [0, T]$, a standard Wiener process $\mathbf{w}_t$, drift coefficient $\boldsymbol{f} \colon \mathbb{R}^d \times [0, T] \to \mathbb{R}^d$, and diffusion coefficient $\boldsymbol{G} \colon [0, T] \to \mathbb{R}^{d \times d}$. Given this forward process, the corresponding reverse-time diffusion process [4, 22] that generates data from noise can be

specified as follows,

$$d\mathbf{x}_t = \left[\boldsymbol{f}(\mathbf{x}_t, t) - \boldsymbol{G}(t)\boldsymbol{G}(t)^\top \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\right] dt + \boldsymbol{G}(t)d\bar{\mathbf{w}}_t, \tag{1}$$

Given an estimate of the *score* $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ of the marginal distribution over $\mathbf{x}_t$ at time $t$, the reverse SDE can then be simulated to recover the original data samples from noise. In practice, the score is intractable to compute and is approximated using a parametric estimator $s_\theta(\mathbf{x}_t, t)$, trained using denoising score matching [2–4, 23]:

$$\min_{\boldsymbol{\theta}} \mathbb{E}_t \mathbb{E}_{p(\mathbf{x}_0)} \mathbb{E}_{p_t(\mathbf{x}_t|\mathbf{x}_0)} \left[\lambda(t)\|\mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0)\|_2^2\right].$$

Above, the time $t$ is usually sampled from a uniform distribution $\mathcal{U}(0, T)$. Given an appropriate choice of $\boldsymbol{f}$ and $\boldsymbol{G}$, the perturbation kernel $p(\mathbf{x}_t|\mathbf{x}_0)$ can frequently be computed analytically (e.g., it is typically Gaussian). Consequently, samples $\mathbf{x}_t$ can be generated in constant time, allowing for fast stochastic gradient updates. The choice of the weighting schedule $\lambda(t)$ plays an essential role during training and can be selected to optimize for likelihood [24] or sample quality [4]. The forward SDE asymptotically converges to an equilibrium distribution (usually a standard isotropic Gaussian) which can be used as a prior to initialize the reverse SDE, which can be simulated using numerical solvers.

## 2.2 A General Recipe for Constructing Stochastic Forward Processes

As has been shown in the MCMC literature [16, 25], it is often beneficial to extend the sampling space into an *augmented* space according to $\mathbf{z} = [\mathbf{x}, \mathbf{m}]^T \in \mathbb{R}^{d_z}$, where $\mathbf{x} \in \mathbb{R}^{d_x}$ is the original state space variable and $\mathbf{m} \in \mathbb{R}^{d_m}$ corresponds to some additional auxiliary dimensions. Simulating the dynamics of the variable $\mathbf{z}$ may have desirable properties, such as faster mixing. Inspired by the naming conventions in statistical physics, we call $\mathbf{x}$ the *position* variable and $\mathbf{m}$ the *momentum* variable. Accordingly, we denote their joint space as *augmented space* (or *phase space* if $\mathbf{x}$ *and* $\mathbf{m}$ have equal dimensions). Note that our notation also captures the scenario where $\mathbf{m}$ is absent (zero-dimensional). We now consider the following form of the stochastic process:

$$d\mathbf{z} = \boldsymbol{f}(\mathbf{z})dt + \sqrt{2\boldsymbol{D}(\mathbf{z})}d\mathbf{w}_t, \tag{2}$$

with drift term $\boldsymbol{f}(\mathbf{z}) \in \mathbb{R}^{d_z}$ and diffusion coefficient $\boldsymbol{D}(\mathbf{z}) \in \mathbb{R}^{d_z \times d_z}$. We assume a desired stationary state distribution $p_s(\mathbf{z})$ specified as

$$p_s(\mathbf{z}) \propto \exp(-H(\mathbf{z})),$$

$$H(\mathbf{z}) = H(\mathbf{x}, \mathbf{m}) = U(\mathbf{x}) + \frac{\mathbf{m}^T M^{-1} \mathbf{m}}{2}, \tag{3}$$

where $H$ represents the Hamiltonian associated with $p_s(\mathbf{z})$. The first term in $H(\mathbf{z})$ represents the potential energy $U(\mathbf{x})$ associated with the configuration $\mathbf{x}$ while the second term represents the kinetic energy associated with the auxiliary (or momentum) variables $\mathbf{m}$ and mass matrix $M\boldsymbol{I}_{d_m}$. In the context of Bayesian inference, [17] propose a framework to elucidate the design space of possible MCMC samplers that sample from $p_s(\mathbf{z})$. In this framework, the drift $\boldsymbol{f}(\mathbf{z})$ can be parameterized as

$$\boldsymbol{f}(\mathbf{z}) = -(\boldsymbol{D}(\mathbf{z}) + \boldsymbol{Q}(\mathbf{z}))\nabla H + \tau(\mathbf{z}), \tag{4}$$

$$\tau_i(\mathbf{z}) = \sum_{j=1}^d \frac{\partial}{\partial \mathbf{z}_j} \left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right),$$

where $\boldsymbol{Q}(\mathbf{z})$ represents a skew-symmetric curl matrix. Furthermore, the following result holds:

**Theorem 2.1** (Yin et. al. [26]). *For the dynamics defined in Eqn. 2, if $\boldsymbol{f}(\mathbf{z})$ is parameterized as in Eqn. 4 with $\boldsymbol{D}(\mathbf{z})$ positive semidefinite and $\boldsymbol{Q}(\mathbf{z})$ skew-symmetric, then the distribution $p_s(\mathbf{z}) \propto \exp\left(-H(\mathbf{z})\right)$ is a stationary distribution for the dynamics.*

Theorem 2.1 implies that for a specific choice of matrices $\boldsymbol{D}(\mathbf{z})$ and $\boldsymbol{Q}(\mathbf{z})$, the process defined in Eqn. 2 always asymptotically samples from the target distribution $p_s(\mathbf{z})$. Moreover, [17] showed in the context of MCMC that the parameterization defined in Eqn. 4 is *complete* as follows:

**Theorem 2.2** (Ma et. al. [17]). *Assume the stochastic process in Eqn. 2 converges to a unique stationary distribution $p_s(\mathbf{z})$. Then, under mild regularity assumptions, there exists a corresponding skew-symmetric matrix $\boldsymbol{Q}(\mathbf{z})$, such that $\boldsymbol{f}(\mathbf{z})$ assumes the form of Eqn. 4.*

We include the proofs for Theorems 2.1 and 2.2 in Appendix A.1 for completeness. These results provide a general recipe for designing forward processes in SGMs.

For the SGM to be a useful forward process, we need it to converge to a simple factorized distribution that serves as the initialization point of the backwards (generative) process. Consequently, we consider the following form of the stationary distribution $p_s(\mathbf{z})$:

$$p_s(\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{0}_{d_x}, \boldsymbol{I}_{d_x}) \mathcal{N}(\mathbf{0}_{d_m}, M\boldsymbol{I}_{d_m}). \tag{5}$$

This form results from setting $U(\mathbf{x}) = \frac{\mathbf{x}^T\mathbf{x}}{2}$ in Eqn. 3. Therefore, for a positive semidefinite matrix $\boldsymbol{D}(\mathbf{z})$ and a skew-symmetric matrix $\boldsymbol{Q}(\mathbf{z})$, the most general class of forward processes which lead to an invariant distribution $p_s(\mathbf{z})$ can be specified by substituting the form of $\nabla H(\mathbf{z})$ (corresponding to $p_s(\mathbf{z})$ defined in Eqn. 5) in Eqn. 4. A similar characterization of forward processes has also been explored in a concurrent work by [27] in the context of likelihood estimation (see Section 5).

## 2.3 Additional constraints on D and Q

Theorems 2.1 and 2.2 show that the proposed forward process parameterization is complete upon specifying the target distribution $p_s(\mathbf{z})$ (such as Eqn. 5). However, we need additional requirements for the resulting generative model and the corresponding training objective to be tractable. Specifically, when using the denoting score matching objective [23], we require the perturbation kernel $p(\mathbf{z}_t|\mathbf{z}_0)$ to be computable in closed form. In practice, this restricts our possible choices for $\boldsymbol{D}(\mathbf{z})$ and $\boldsymbol{Q}(\mathbf{z})$ to constant matrices (i.e., independent of the state variable $\mathbf{z}$). Yet, even with this requirement, the framework provides a large design space of models. We provide several examples of existing SGMs that can be understood as special cases of our recipe in Appendix A.2. We stress that training paradigms other than denoting score matching (e.g., such as Sliced Score matching [28]) may enable a wider range of possible models with non-constant matrices $D$ and $Q$.

## 3 Phase Space Langevin Diffusion

We next use the proposed recipe to construct a specific SGM with favorable properties.

### 3.1 Model Definition

We restrict the family of forward processes considered in this work by constraining $\boldsymbol{D}(\mathbf{z})$ and $\boldsymbol{Q}(\mathbf{z})$ as constant matrices, i.e., independent of state $\mathbf{z}$. Moreover, we assume that $\mathbf{x}$ and $\mathbf{m}$ have the same dimension d, i.e. $\mathbf{z} \in \mathbb{R}^{2d}$. Consequently, the drift $f(\mathbf{z})$ becomes affine in $\mathbf{z}$ and the perturbation kernel $p(\mathbf{z}_t|\mathbf{z}_0)$ can be computed analytically [29]. Among the possible samplers, we choose a specific form involving $d-$dimensional position and momentum coordinates, $\mathbf{z}_t = [\mathbf{x}_t, \mathbf{m}_t]^T$ where $\mathbf{x}_t \in \mathbb{R}^d$, $\mathbf{m}_t \in \mathbb{R}^d$. Our choice for $\boldsymbol{D}(\mathbf{z})$ and $\boldsymbol{Q}(\mathbf{z})$ is as follows:

$$\boldsymbol{D} := \frac{\beta}{2}\left(\begin{pmatrix} \Gamma & 0 \\ 0 & M\nu \end{pmatrix} \otimes \boldsymbol{I}_d\right), \quad \boldsymbol{Q} := \frac{\beta}{2}\left(\begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \otimes \boldsymbol{I}_d\right). \tag{6}$$

Above, $\Gamma$, $M$, $\nu$ and $\beta$ are positive scalars. Along with these choices of $\boldsymbol{D}$ and $\boldsymbol{Q}$, we have $\tau(\mathbf{z}) = \mathbf{0}$. The resulting forward process is given by:

$$d\mathbf{z}_t = \boldsymbol{f}(\mathbf{z}_t)dt + \boldsymbol{G}(t)d\mathbf{w}_t, \tag{7}$$

$$\boldsymbol{f}(\mathbf{z}_t) = \left(\frac{\beta}{2}\begin{pmatrix} -\Gamma & M^{-1} \\ -1 & -\nu \end{pmatrix} \otimes \boldsymbol{I}_d\right)\mathbf{z}_t, \quad \boldsymbol{G}(t) = \sqrt{2D(\mathbf{z}_t)} = \begin{pmatrix} \sqrt{\Gamma\beta} & 0 \\ 0 & \sqrt{M\nu\beta} \end{pmatrix} \otimes \boldsymbol{I}_d.$$

We denote the form of the SDE in Eqn. 7 as the *Phase Space Langevin Diffusion (PSLD)*. Note that PSLD generalizes Critically Damped Langevin Diffusion (Critically Damped Langevin Diffusion (CLD)) proposed in [14], which can be obtained by setting $\Gamma = 0$, $\bar{\nu} = M\nu$, and $\bar{\beta} = \frac{\beta}{2}$. Like CLD, the parameter $M^{-1}$ couples the data space state $\mathbf{x}_t$ with the auxiliary state $\mathbf{m}_t$. The parameters $\beta$, $\Gamma$, and $\nu$ control the amount of noise in the forward SDE. Without loss of generality, we use a time-independent $\beta$. However, unlike CLD or any physical system, PSLD adds stochastic noise in the data space *in addition* to the noise injected into the momentum component of phase space. While we are not aware of any physical system that displays such behavior, it is a valid stochastic process

compatible with our framework. Our experiments reveal the strong benefits of having these two independent noise sources.

Furthermore, CLD [14] proposes setting $\bar{\nu}^2 = 4M$, corresponding to critical damping in a physical system. Under critical damping, an ideal balance is achieved between the oscillatory Hamiltonian dynamics and the noise-injecting Ohrnstein-Uhlenbeck (OU) term, leading to faster convergence to equilibrium. We generalize this line of argument in Appendix B.1, where we derive $(\Gamma - \nu)^2 = 4M^{-1}$ as the equivalent condition for critical damping in PSLD. Throughout this work, we choose $\Gamma$, $\nu$, and $M^{-1}$ such that the critical damping condition in PSLD is satisfied.

## 3.2 PSLD Training

Since the drift coefficient in PSLD is affine, the perturbation kernel $p(\mathbf{z}_t|\mathbf{z}_0)$ of PSLD can be computed analytically. We can then use DSM to learn the score function $s_\theta(\mathbf{z}_t, t)$. More specifically, following the derivation in [24], it can be shown that the Maximum-Likelihood (ML) based DSM objective for PSLD can be specified as (Proof in Appendix B.2.1)

$$\min_{\boldsymbol{\theta}} \mathbb{E}_t \mathbb{E}_{p(\mathbf{z}_0)} \mathbb{E}_{p_t(\mathbf{z}_t|\mathbf{z}_0)} \Big[ \Gamma \beta \mathcal{L}_x(\theta, \mathbf{z}_t, \mathbf{z}_0) + M\nu\beta \mathcal{L}_m(\theta, \mathbf{z}_t, \mathbf{z}_0) \Big], \tag{8}$$

$$\mathcal{L}_x = \|s_{\boldsymbol{\theta}}(\mathbf{z}_t, t)|_{0:d} - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{z}_t|\mathbf{z}_0)\|_2^2, \quad \mathcal{L}_m = \|s_{\boldsymbol{\theta}}(\mathbf{z}_t, t)|_{d:2d} - \nabla_{\mathbf{m}_t} \log p_t(\mathbf{z}_t|\mathbf{z}_0)\|_2^2, \tag{9}$$

where $s_{\boldsymbol{\theta}}(\mathbf{z}_t, t)|_{0:d}$ and $s_{\boldsymbol{\theta}}(\mathbf{z}_t, t)|_{d:2d}$ represent the first and the last $d$ components of the vector $s_{\boldsymbol{\theta}}(\mathbf{z}_t, t)$ respectively. In the above DSM objective, the perturbation kernel $p(\mathbf{z}_t|\mathbf{z}_0) = \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$ is a multivariate Gaussian while $p(\mathbf{z}_0) = p(\mathbf{x}_0)\mathcal{N}(\mathbf{m}_0; 0, M\gamma \mathbf{I}_d)$, where $p(\mathbf{x}_0)$ is the data distribution. In this work, we reformulate the DSM objective in Eqn. 8 as follows (also see Appendix B.2.1):

$$\min_{\boldsymbol{\theta}} \mathbb{E}_t \mathbb{E}_{p(\mathbf{z}_0)} \mathbb{E}_{p_t(\mathbf{z}_t|\mathbf{z}_0)} \Big[ \lambda(t) \|s_{\boldsymbol{\theta}}(\mathbf{z}_t, t) - \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t|\mathbf{z}_0)\|_2^2 \Big].$$

Furthermore, due to its gradient variance reduction properties, we instead use the Hybrid Score Matching (HSM) objective [14] by marginalizing out the momentum variables $\mathbf{m}_0$ as $p(\mathbf{z}_t|\mathbf{x}_0) = \int p(\mathbf{z}_t|\mathbf{x}_0, \mathbf{m}_0)p(\mathbf{m}_0)d\mathbf{m}_0$. Since both distributions $p(\mathbf{z}_t|\mathbf{x}_0, \mathbf{m}_0)$ and $p(\mathbf{m}_0)$ are Gaussian, $p(\mathbf{z}_t|\mathbf{x}_0)$ will also be a Gaussian.

**Score Network Parameterization:** Since the perturbation kernel $p(\mathbf{z}_t|\mathbf{x}_0)$ in the HSM objective is also a multivariate Gaussian, we have $p(\mathbf{z}_t|\mathbf{x}_0) \sim \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$. Furthermore, let $\boldsymbol{\Sigma}_t = \boldsymbol{L}_t \boldsymbol{L}_t^T$ be the Cholesky factorization of the matrix $\boldsymbol{\Sigma}_t$. We have

$$\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{x}_0) = -\boldsymbol{\Sigma}_t^{-1}(\mathbf{z}_t - \boldsymbol{\mu}_t) = -\boldsymbol{L}_t^{-T} \boldsymbol{\epsilon}, \tag{10}$$

where $\boldsymbol{L}_t^{-T}$ is the transposed inverse of the $\boldsymbol{L}_t$ and $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}_{2d}, \mathbf{I}_{2d})$. Therefore, we parameterize our score function estimator as $s_{\boldsymbol{\theta}}(\boldsymbol{z}_t, t) = -\boldsymbol{L}_t^{-T}\boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\mathbf{z}_t, t)$. Although alternative parameterizations of the score network $s_\theta(\mathbf{z}_t, t)$ like *mixed score* can be possible [14, 30, 31], we do **not** explore such parameterizations in this work and leave further exploration to future work. We provide additional details on the score network parameterization in PSLD in Appendix B.2.2 and the analytical form of the perturbation kernel $p(\mathbf{z}_t|\mathbf{x}_0)$ in Appendix B.3.

**Final Training Objective:** Using our score parameterization from Eqn. 10 with $\lambda(t) = \frac{1}{\|\boldsymbol{L}_t^{-T}\|_2^2}$, we get the following *epsilon-prediction* form of the HSM objective (See Appendix B.2.3 for a complete derivation):

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{p(\mathbf{x}_0)} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{I}_d)} \big[ \|\boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{\mu}_t + \boldsymbol{L}_t \boldsymbol{\epsilon}, t) - \boldsymbol{\epsilon}\|_2^2 \big]. \tag{11}$$

The epsilon-prediction objective has been shown to generate superior sample quality [3, 4, 14]. In this work, we optimize for sample quality and therefore use this objective for training all models. One key difference between the objective in Eqn. 11 and the HSM objective in CLD is that, unlike CLD, we predict the full 2d-dimensional $\boldsymbol{\epsilon}$ due to the structure of our diffusion coefficient $\boldsymbol{G}(t)$ (see Appendix B.2 for more details). Therefore, for a non-zero $\Gamma$, the neural-net-based score predictor in PSLD has twice the number of output channels as in CLD. However, the increase in parameters due to this architectural update is negligible.

### 3.3 PSLD Sampling

Following the result from [4], the reverse process SDE corresponding to the forward process SDE defined in Eqn. 8 can be formulated as follows:

$$d\bar{\mathbf{z}}_t = \bar{\boldsymbol{f}}(\bar{\mathbf{z}}_t)dt + \boldsymbol{G}(T-t)d\bar{\mathbf{w}}_t \tag{12}$$

$$\bar{\boldsymbol{f}}(\bar{\mathbf{z}}_t) = \frac{\beta}{2}\begin{pmatrix} \Gamma\bar{\mathbf{x}}_t - M^{-1}\bar{\mathbf{m}}_t + 2\Gamma\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{0:d} \\ \bar{\mathbf{x}}_t + \nu\bar{\mathbf{m}}_t + 2M\nu\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{d:2d} \end{pmatrix}, \quad \boldsymbol{G}(T-t) = \begin{pmatrix} \sqrt{\Gamma\beta} & 0 \\ 0 & \sqrt{M\nu\beta} \end{pmatrix} \otimes \boldsymbol{I}_d \tag{13}$$

where $\bar{\mathbf{z}}_t = \mathbf{z}_{T-t}$, $\bar{\mathbf{x}}_t = \mathbf{x}_{T-t}$, $\bar{\mathbf{m}}_t = \mathbf{m}_{T-t}$. We can simulate this reverse process SDE using standard numerical SDE solvers like the Euler-Maruyama (EM) sampler [32]. As an alternative, Dockhorn et al. [14] propose SSCS: a symmetric splitting-based integrator and show that SSCS exhibits a better speed-sample quality tradeoff than EM. Consequently, we extend SSCS for PSLD by using the following splitting formulation:

$$\begin{pmatrix} d\bar{\mathbf{x}}_t \\ d\bar{\mathbf{m}}_t \end{pmatrix} = \underbrace{\frac{\beta}{2}\begin{pmatrix} -\Gamma\bar{\mathbf{x}}_t - M^{-1}\bar{\mathbf{m}}_t \\ \bar{\mathbf{x}}_t - \nu\bar{\mathbf{m}}_t \end{pmatrix} dt + \boldsymbol{G}(T-t)d\bar{\mathbf{w}}_t}_{\text{Analytical-term}} + \underbrace{\beta\begin{pmatrix} \Gamma\bar{\mathbf{x}}_t + \Gamma\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{0:d} \\ \nu\bar{\mathbf{m}}_t + M\nu\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{d:2d} \end{pmatrix} dt}_{\text{Score-term}} \tag{14}$$

Indeed for $\Gamma = 0$, the sampler in Eqn. 14 resembles the SSCS sampler proposed in [14]. It is worth noting that despite an updated formulation, the order of the SSCS sampler, as analyzed in [14], remains unchanged. We discuss the exact solution of the analytical part of the Modified-SSCS sampler in Eqn. 14 and other relevant details in Appendix B.4.2.

## 4 Experiments

**Datasets**: We run experiments on three datasets: CIFAR-10 [18], CelebA [19] at 64 x 64 resolution and the AFHQv2 [33] dataset at 128 x 128 resolution.

**Baselines**: We primarily compare PSLD with two popular SGM baselines: VP-SDE [4] and CLD [14] (a particular case of PSLD with $\Gamma = 0$). For PSLD and CLD, unless specified otherwise, we operate in the critical damping regime with a fixed $M^{-1} = 4$ and therefore choose $\Gamma$ and $\nu$ accordingly ($\nu = 2\sqrt{M^{-1}} + \Gamma$, $\nu \geq 0$, $\Gamma \geq 0$).

**Metrics**: We use the FID [20] score for quantitatively assessing sample quality, while we use NFE (Number of Function Evaluations) to assess the sampling efficiency of all methods.

We provide full implementation details in Appendix C. The rest of our experimental section is organized as follows: Firstly, we compare the state-of-the-art performance of PSLD with popular SGM baselines on unconditional image generation. We show that PSLD outperforms competing baselines for similar compute budgets. Secondly, as an ablation experiment, we empirically and theoretically analyze the impact of the SDE parameters $\Gamma$ and $\nu$ on downstream sample quality in PSLD. Furthermore, we analyze the speed-quality trade-off in PSLD and show that PSLD yields better sample quality than competing baselines across four different sampler settings. Lastly, we show that pre-trained unconditional PSLD models can be used for downstream tasks like class-conditional image synthesis and image inpainting.

### 4.1 State-of-the-art Comparisons

**Setup**: We now compare the sample quality of our proposed method with existing popular SGM methods on the CIFAR-10 and CelebA-64 datasets for unconditional image synthesis. We use PSLD with $\Gamma \in \{0.01, 0.02\}$ for CIFAR-10 and PSLD with $\Gamma = 0.005$ for CelebA-64 for state-of-the-art (SOTA) comparisons (See Section 4.2 for a theoretical and empirical justification of these choices). Moreover, we use the training objective in Eqn. 11 **without** any alternative score parameterizations (like *mixed score*[14, 30]) to train our models for SOTA comparisons. Unless specified otherwise, we perform sampling using the EM sampler with Uniform Striding (US) for CIFAR-10 and Quadratic Striding (QS) for CelebA-64 for the SDE setup and report FID scores on 50k samples (denoted as FID@50k). We include full details on our SDE and ODE solver setup for SOTA analysis in Appendix C.5. We report the FID scores for most competing methods for a maximum sampling

| Method | Size | NFE | FID@50k ($\downarrow$) |
|---|---|---|---|
| **Ours (Baseline)** | | | |
| CLD (w/o MS) | 97M | 1000 | 2.41 |
| **Ours (Proposed)** | | | |
| PSLD ($\Gamma$=0.02) | 39M | 1000 | 2.80 |
| PSLD ($\Gamma$=0.01) | 55M | 1000 | 2.34 |
| PSLD ($\Gamma$=0.02) | 55M | 1000 | 2.30 |
| PSLD ($\Gamma$=0.01) | 97M | 1000 | 2.23 |
| PSLD ($\Gamma$=0.02) | 97M | 1000 | **2.21** |
| CLD (w/ MS) [14] | 108M | 1000 | 2.27 |
| VPSDE (deep)[†] [4] | 108M | 1000 | 2.46 |
| VESDE (deep)[†] [4] | 108M | 1000 | 2.43 |
| DDPM [3] | 35.7M | 1000 | 3.17 |
| iDDPM [34] | - | 1000 | 2.90 |
| DiffuseVAE [35] | 35.7M | 1000 | 2.80 |
| NCSNv2 [36] | - | - | 10.87 |
| NCSN [2] | - | 1000 | 25.32 |
| VDM [37] | - | 1000 | 7.41 |

Table 1: PSLD (SDE) sample quality comparisons for CIFAR-10. PSLD outperforms competing SDE baselines for a similar sampling budget. FID computed using 50k samples. MS: Mixed Score †: Results from [14].

| Method | Size | NFE | FID@50k ($\downarrow$) |
|---|---|---|---|
| **Ours (Baseline)** | | | |
| CLD (w/o MS) | 97M | 352 | 2.80 |
| **Ours (Proposed)** | | | |
| PSLD ($\Gamma$=0.01) | 55M | 243 | 2.41 |
| PSLD ($\Gamma$=0.02) | 55M | 232 | 2.40 |
| PSLD ($\Gamma$=0.01) | 97M | 246 | **2.10** |
| PSLD ($\Gamma$=0.02) | 97M | 231 | 2.31 |
| **Ours (Proposed)** | | | |
| PSLD ($\Gamma$=0.01) | 97M | 159 | 2.13 |
| PSLD ($\Gamma$=0.02) | 97M | 159 | 2.34 |
| LSGM[†] [30] | 100M | 131 | 4.60 |
| LSGM [30] | 476M | 138 | **2.10** |
| VPSDE[†] [4] | 108M | 141 | 2.76 |
| CLD (w/ MS) [14] | 108M | 147 | 2.71 |
| CLD (w/ MS) [14] | 108M | 312 | 2.25 |
| ScoreFlow (VP) [24] | 108M | - | 5.34 |
| Flow Matching (w/ OT) [38] | - | 142 | 6.35 |
| DDIM (VPSDE)[†] [39] | 108M | 150 | 3.15 |

Table 2: PSLD (ODE) sample quality comparisons for CIFAR-10. PSLD outperforms most competing ODE baselines. FID computed using 50k samples. MS: Mixed Score. †: From [14].

budget of N=1000 while reporting model sizes whenever available for CIFAR-10 for fair comparisons.

**Main Observations**: Table 1 compares CIFAR-10 sample quality between different methods using stochastic sampling. Our proposed method with $\Gamma = 0.02$ and 39M parameters achieves an FID score of 2.80, outperforming the DDPM [3] baseline while performing comparably with DiffuseVAE [35], for similar model sizes. It is worth noting that DiffuseVAE refines samples generated from a VAE [40] using a DDPM backbone and is complementary to our work. *Furthermore, our larger PSLD model achieves an FID of 2.21, which is better than CLD [14] (with or without the Mixed Score (MS) parameterization) and (VP/VE)-SDE baselines for similar NFE budget and model sizes.* For CelebA-64 (Table 3), PSLD outperforms the VP/VE-SDE baselines by a significant margin while requiring only 250 NFEs.

We next analyze ODE sample quality in PSLD. Table 2 compares CIFAR-10 sample quality between different methods using ODE-based samplers. PSLD with $\Gamma = 0.01$ achieves an FID score of 2.10 and outperforms most competing methods except LSGM [30]. Though the original LSGM model is more than four times the size of our SOTA model, PSLD performs comparably with LSGM. When scaled to a similar size, LSGM performs much worse than PSLD (FID: 4.60 for LSGM-100M vs. 2.10 for PSLD). We note that EDM [31] achieves an FID of 2.05 on unconditional CIFAR-10 generation (without data augmentation) by analyzing several design choices associated with diffusion models (like score network architectures, loss preconditioning, and sampler design). We did not explore this line of research but note that their approach complements our proposed method, and exploring some of these design choices in the context of PSLD can be an exciting future direction.

Interestingly, PSLD with the ODE setup obtains a better FID score than the SDE setup (FID: 2.21 for SDE vs. 2.10 for ODE) while requiring around four times lesser NFEs. Moreover, when using a solver tolerance of $1e^{-4}$, PSLD achieves an FID score of 2.13, comparable to the best FID of 2.10 while reducing NFEs significantly. *This tradeoff is worse for other SGMs like CLD and VP-SDE (Table 2).* We report additional SOTA results in Appendix D.3.

### 4.2 Impact of $\Gamma$ and $\nu$ on PSLD Sample Quality

**Setup and Baselines**: Since adding stochastic noise in both the data and the momentum space is one of the primary aspects of PSLD, we now analyze the impact of the choice of $\Gamma$ and $\nu$ on downstream sample quality. For subsequent experimental results, we use our smaller ablation models (for PSLD and relevant baselines) for comparisons. Table 4 shows the impact of varying $\Gamma$ on sample quality

| Method | NFE | FID@50k |
|---|---|---|
| **(Ours) PSLD ($\Gamma = 0.005$)** | 250 | **2.01** |
| **(Ours) PSLD ($\Gamma = 0.005$, ODE)** | 244 | 2.56 |
| PNDM [41] | 250 | 2.71 |
| DDIM [39] | 250 | 4.44 |
| VPSDE [4] | 1000 | 2.32 |
| Gamma DDPM [42] | 1000 | 2.92 |
| DDPM [3] | 1000 | 3.26 |
| DiffuseVAE [35] | 1000 | 3.97 |
| VESDE [4] | 2000 | 3.95 |
| NCSN (w/ denoising) [2] | - | 25.3 |
| NCSNv2 (w/ denoising) [36] | - | 10.23 |

Table 3: PSLD sample quality comparisons for CelebA-64. FID computed using 50k samples.

| | CIFAR-10 (39M) | | CelebA-64 (66M) | |
|---|---|---|---|---|
| $\Gamma$ | FID@50k (EM-QS) | FID@50k (EM-US) | FID@10k (EM-QS) | FID@10k (EM-US) |
| 0 | 3.64 | 3.60 | 4.59 | 4.60 |
| 0.005 | 3.42 | 3.34 | **4.17** | 4.37 |
| 0.01 | 3.15 | 2.94 | 4.22 | 4.34 |
| 0.02 | 3.26 | **2.80** | 4.43 | 4.52 |
| 0.25 | 4.99 | 9.48 | 93.99 | 95.13 |

Table 4: Impact of increasing $\Gamma$ (with fixed $M^{-1}$) on sample quality (NFE=1000). FID computed using 50k and 10k samples for the CIFAR-10 and CelebA-64 datasets, respectively. QS: Quadratic Striding, US: Uniform Striding.



Figure 2: Impact of increasing $\Gamma = \{0, 0.005, 0.02, 0.25\}$ (Top to Bottom) on CelebA-64 sample quality. The best sample quality is achieved at $\Gamma = 0.005$ (Second Row) while increasing $\Gamma$ to 0.25 results in loss of high-frequency image features.

for the CIFAR-10 and CelebA-64 datasets. Our ablation CLD baseline (PSLD with $\Gamma = 0, \nu = 4$) achieves an FID of 3.60 using the Euler-Maruyama (EM) sampler with Uniform striding (US) and 3.64 using the EM sampler with Quadratic striding (QS). Our results are comparable with the FID of 3.56 obtained by [14] for their CLD ablation model on CIFAR-10 without using the mixed-score parameterization. Our VP-SDE ablation baseline (not shown in Table 4) obtains an FID of 3.19 using EM-US (with 1000 NFEs).

**Main Observations**: We observe that *setting $\Gamma$ to a non-zero value within a specific range improves sample quality significantly over CLD*. Specifically, our ablation CIFAR-10 model achieves FID scores of 2.94 and 2.80 for $\Gamma = 0.01$ and $\Gamma = 0.02$ respectively (with EM-US) and outperforms our VP-SDE and CLD baselines without using alternative score-network parameterizations like mixed-score which is crucial for competitive performance of CLD [14]. We make a similar observation for the CelebA-64 dataset on which our model achieves the best FID of 4.17 using EM-QS and outperforms our CLD baseline (FID: 4.59). Interestingly, the sample quality worsens for both datasets on increasing $\Gamma$ outside a range. For instance, for CIFAR-10, further increasing $\Gamma$ from 0.02 to 0.04 (not shown in Table 4) resulted in an increase in FID from $2.80$ to $2.95$. Consequently, the sample quality for both datasets is the worst at $\Gamma = 4.25, \nu = 0.25$. We also note that EM-US works better than EM-QS for CIFAR-10 and vice-versa for CelebA-64.

Figure 2 further validates our findings qualitatively for the CelebA-64 dataset where for $\Gamma = 0.25$, the score network can only recover high-level semantic structures (like gender and glasses, among others) but is unable to recover high-frequency details. Since the diffusion denoiser recovers most high-frequency information in the low-timestep regime, these observations suggest denoising issues near low-timestep indices. We next provide a formal justification for this observation.

| | NFE (FID@10k ↓) | | | | |
|---|---|---|---|---|---|
| Sampler | Method | 50 | 100 | 250 | 500 | 1000 |
| EM-QS | CLD | 25.01 | 8.91 | 5.97 | 5.61 | 5.7 |
| | VP-SDE | **17.72** | 7.45 | 5.59 | 5.51 | 5.51 |
| | (Ours) PSLD | 19.94 | **7.33** | **5.26** | **5.20** | **5.28** |
| EM-US | CLD | 119.68 | 45.60 | **9.08** | 5.71 | 5.65 |
| | VP-SDE | **84.54** | 41.93 | 12.61 | 5.92 | 5.19 |
| | (Ours) PSLD | 100.62 | **39.96** | 11.26 | **5.45** | **4.82** |
| SSCS-QS | CLD | 21.31 | 8.37 | 5.82 | 5.75 | 5.69 |
| | (Ours) PSLD | **16.12** | **7.16** | **5.36** | **5.35** | **5.27** |
| SSCS-US | CLD | 75.45 | 24.74 | 6.09 | 5.74 | 5.78 |
| | (Ours) PSLD | **72.42** | **20.46** | **5.19** | **4.92** | **5.29** |

Table 5: PSLD exhibits better speed vs. sample quality tradeoffs over competing baseline SDEs (CLD and VP-SDE) on CIFAR-10 across four samplers configurations. The rightmost five columns indicate NFEs, with **bold** indicating the best result for that sampler. QS: Quadratic Striding, US: Uniform Striding. See Appendix D.2 for extended results.

| Method | $\log_{10}$ tol | FID@10k ($\downarrow$) | Avg. NFE |
|---|---|---|---|
| CLD (Baseline) | -5 | 5.54 | 280 |
| | -4 | 5.62 | 196 |
| | -3 | 6.54 | 147 |
| | -2 | 9.98 | 86 |
| | -1 | 397.1 | 27 |
| PSLD ($\Gamma = 0.02$) | -5 | **4.79** | 228 |
| | -4 | 4.84 | 158 |
| | -3 | 5.09 | 111 |
| | -2 | 16.11 | 69 |
| | -1 | 418.779 | 27 |
| VPSDE | -5 | 5.91 | 123 |

Table 6: PSLD exhibits better speed vs. sample quality tradeoffs over competing baselines on CIFAR-10 using a black-box ODE solver. $\log_{10}$ tol indicates the ODE sampler (RK45) tolerance. **Bold** indicates best result for that column.

**Theoretical justification of adding stochasticity in the position space**: Since PSLD involves adding stochasticity in both the data and the momentum space, during training, we need to predict the noise $\epsilon_\theta^x(\mathbf{z}_t, t)$ and $\epsilon_\theta^m(\mathbf{z}_t, t)$ in both the data and the momentum space respectively. Therefore, it is unclear why PSLD leads to better sample quality than CLD since predicting both noise components can lead to additional sources of errors during sampling.

However, in the context of the EM sampler, we find (see Appendix D.1) that setting a *small* non-zero $\Gamma$ can significantly suppress prediction errors from $\epsilon_\theta^x(\mathbf{z}_t, t)$ at the expense of introducing negligible extra errors from $\epsilon_\theta^m(\mathbf{z}_t, t)$. Contrarily, using larger values of $\Gamma$ results in scaling the prediction errors from $\epsilon_\theta^x(\mathbf{z}_t, t)$ by a significant factor, especially in the low-timestep regime, leading to worse sample quality with significant degradations in high-frequency sample details as observed in Figure 2.

Therefore, intuitively, $\Gamma$ *introduces a trade-off between error contribution from both noise predictors* $\epsilon_\theta^x(\mathbf{z}_t, t)$ *and* $\epsilon_\theta^m(\mathbf{z}_t, t)$ *with small values of* $\Gamma$ *providing a favorable trade-off which improve overall sample quality*. As a general guideline, we find $\Gamma = 0.01$ to work well across datasets. Figure 1 shows some qualitative samples generated from PSLD trained on the AFHQv2 [33] dataset with $\Gamma = 0.01, \nu = 4.01$.

## 4.3 Sample Speed vs. Quality Tradeoffs for PSLD

**Sampler Setup**: Since the tradeoff between sample quality and the number of reverse sampling steps required is crucial for any SGM backbone, we now examine this tradeoff for PSLD for the CIFAR-10 dataset (See Appendix D.2 for extended results on the CelebA-64 dataset). We use our VP-SDE and PSLD with $\Gamma = 0$ (corresponding to CLD) ablation models as comparison baselines. Furthermore, we use combinations of the EM and SSCS samplers with Uniform (US) and Quadratic (QS) timestep striding as different sampler settings to benchmark the performance of all methods. It is worth noting that the SSCS sampler can only be used for augmented SGMs like CLD and PSLD. For ODE-based comparisons, we use the probability flow ODE setup with RK45 [43] solver (see Appendix C.5 for more details). Lastly, we measure sample quality using FID computed for 10k samples.

**Main Observations**: Table 5 shows a comparison between FID scores for our best performing PSLD models (corresponding to $\Gamma = 0.01$ and $\Gamma = 0.02$) and our VP-SDE and CLD baselines from Section 4.2 for the CIFAR-10 dataset across $N \in \{50, 100, 250, 500, 1000\}$ steps. We primarily observe that *PSLD outperforms the VP-SDE and CLD baselines across all comparison points* with the most significant differences at lower NFE (Network Function Evaluations) values. For instance, PSLD ($\Gamma = 0.02$) achieves the best FID of 16.12 at NFE=50 compared to 17.72 and 21.31 by VP-SDE and CLD, respectively. Moreover, for most sampler settings across methods, SSCS performs better in the low NFE regime ($N \leq 250$), while EM performs better for a higher number of NFEs. A

Figure 3: (Left) Class-conditional results on CIFAR-10: *Truck*, *Airplane*, and *Automobile* from Top to Bottom (two rows each). (Middle) Class conditional results on AFHQv2: *Dogs*, *Cats* and *Others* from Top to Bottom (one row each). (Right) Inpainting results on AFHQv2. The columns represent the original, corrupted, and imputed samples, respectively, from left to right.

| Method | FID (Train) | FID (Test) |
|---|---|---|
| CLD | 1.01 | 7.10 |
| (Ours) PSLD ($\Gamma = 0.01$) | **0.85** | **6.93** |

Table 7: PSLD outperforms CLD on Image inpainting for the AFHQv2 dataset. FID (lower is better) is computed on the full train and test sets.

similar observation was made in [14]. Similarly, Quadratic striding works much better in the low NFE regime, while Uniform striding works better when using a higher number of NFEs ($N > 500$).

We next compare our best-performing ablation model (PSLD with $\Gamma = 0.02$) with our VP-SDE and CLD baselines using the probability flow ODE setup across multiple tolerance levels on the CIFAR-10 dataset. Table 6 compares FID scores (computed on 10k samples) for all methods. Like our SDE setup, *PSLD ($\Gamma=0.02$) outperforms both baselines in sample quality for similar NFE budgets*. Moreover, across the same solver tolerance level, PSLD requires fewer NFEs on average than its CLD counterpart while yielding better sample quality. Lastly, we found using black-box solvers to further improve sample quality compared to the SDE baseline at a tolerance level of 1e-5 for both our PSLD and CLD models for CIFAR-10 (FID@10k=4.97 for PSLD for ODE vs. FID@10k=4.84 for EM-QS (N=1000) with $\Gamma = 0.02$). This observation is consistent with the ODE comparison results presented in Section 4.1.

### 4.4 Conditional Generation with PSLD

Following prior work [4, 5], given some conditioning information $\mathbf{y}$, an unconditional pre-trained score network $\boldsymbol{s}_\theta(\mathbf{z}_t, t)$ can be used for sampling from the distribution $p(\mathbf{z}_t|\mathbf{y})$ in PSLD. More specifically,

$$\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{y}) = \nabla_{\mathbf{z}_t} \log p(\mathbf{y}|\mathbf{z}_t) + \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \tag{15}$$

$$\approx \nabla_{\mathbf{z}_t} \log p(\mathbf{y}|\mathbf{z}_t) + \boldsymbol{s}_\theta(\mathbf{z}_t, t) \tag{16}$$

We can then use the estimate of $\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{y})$ in Eqn. 16 to sample from the following SDE for conditional generation:

$$d\mathbf{z}_t = \left[ \boldsymbol{f}(\mathbf{z}_t) - \boldsymbol{G}(t)\boldsymbol{G}(t)^T \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{y}) \right] dt + \boldsymbol{G}(t)d\mathbf{w}_t. \tag{17}$$

Figure 3 illustrates class conditional samples for the CIFAR-10 and the AFHQv2 datasets obtained by training an additional *time-dependent classifier* $p(\mathbf{y}|\mathbf{z}_t)$ to compute $\nabla_{\mathbf{z}_t} \log p(\mathbf{y}|\mathbf{z}_t)$, followed by sampling from the SDE in Eqn. 17 (full implementation details in Appendix D.4). Similarly, we can perform data imputation by setting the conditioning signal $\mathbf{y} = \bar{\mathbf{z}}_0$ where $\bar{\mathbf{z}}_0$ is the observed part of the input data $\mathbf{z}_0$ (See Figure 3). For image inpainting, PSLD exhibits a better perceptual quality of inpainted samples over CLD on the AFHQv2 dataset (See Table 7). We include a complete derivation for inpainting and an analogous framework to [4] for solving inverse problems using PSLD with additional conditional synthesis results in Appendix D.4.

# 5 Related Work

**Advances in Diffusion Models**: Following the seminal work on diffusion (a.k.a score-based) models [1–4], there has been much recent progress in advancing unconditional [5, 7, 14, 30, 34, 44–46] and conditional [4, 7, 35, 47, 48] diffusion models for a variety of downstream tasks like text-to-image synthesis [49, 50], image super-resolution [47, 51] and video generation [11, 52–54]. Our work is closely related to CLD [14], which is motivated by Langevin heat baths in statistical mechanics [55]. However, our method is not directly motivated by physical interpretation but rather directly constructed from our proposed drift parameterization. Another line of research in SGMs is to perform score-based modeling in the latent space [7, 30, 56] of a powerful autoencoder [57, 58]. Such approaches have been shown to improve the sampling time in SGMs. Therefore, since we propose a novel diffusion model backbone, most existing advances in diffusion models complement PSLD.

**Sampler Design in Diffusion Models**: Improving the speed-vs-quality tradeoff in SGMs is a fundamental area in diffusion model research [39, 41, 59–62]. One popular approach to speed up diffusion model sampling is DDIM [39]. [61] show that DDIM can be cast as an exponential integrator and propose further improvements. [62] further leverage these improvements to propose a *generalized-DDIM* (gDDIM) method for CLD. It is worth noting that gDDIM parameterization requires predicting the score w.r.t both the data and the auxiliary variables and is directly compatible with PSLD. Another line of research involves training to speed up diffusion sampling. GENIE [63] proposes to utilize higher-order Taylor methods during training to speed up DDIM sampling. Alternatively, distillation-based approaches distill a teacher into a student diffusion model progressively [46, 64] or otherwise [65]. Therefore, exploring some of these directions in the context of PSLD would be interesting.

**Auxiliary Diffusion Models**: In a concurrent work, [27] define an ELBO for multivariate diffusion models (MDM) and introduce a similar recipe as ours to design new diffusion processes. While [27] optimize for likelihood estimates, we primarily focus on sample quality in this work. Both works illustrate a different perspective on the advantages of constructing a generic recipe for designing diffusion processes and, therefore, complementary. Another recent work, Flexible Diffusion [66], exploits the geometry of the data manifold to parameterize the forward process. The proposed framework is complete under linear drift. However, our parameterization makes no such assumptions.

# 6 Conclusion

We presented a recipe for constructing forward process parameterization for diffusion processes that guarantees convergence to a prespecified stationary distribution, such as a Gaussian. We use the proposed recipe to construct a novel diffusion process: Phase Space Langevin Diffusion(PSLD) which achieves excellent sample quality with better speed-vs-quality tradeoffs compared to existing baselines like the VP-SDE and CLD on standard image-synthesis benchmarks. We left the exploration of potentially performance-improving design choices such as alternative score network parameterizations and loss weighting [31] as directions for future work.

While this work only explores stochastic samplers with a single auxiliary "momentum" variable $\mathbf{m}$ (of the same dimension as $\mathbf{x}$), exploring other design choices of $\boldsymbol{D}(\mathbf{z})$ and $\boldsymbol{Q}(\mathbf{z})$ [27], which lead to higher-order stochastic samplers (like the Nosé-Hoover Thermostat) could also be an interesting research direction. Furthermore, our current choices of $\boldsymbol{D}(\mathbf{z})$ and $\boldsymbol{Q}(\mathbf{z})$ are limited to constant matrices due to relying on denoising score matching. Therefore, the proposed parameterization offers a complementary framework for designing diffusion generative models trained using alternative score-matching techniques.

Lastly, our proposed recipe is only *complete* under the assumption that both $\mathbf{x}$ and $\mathbf{m}$ are required to converge to prescribed marginals $p(\mathbf{x})$ and $p(\mathbf{m})$. Without this requirement on $\mathbf{m}$, the design space of samplers is potentially larger (as has been pointed out in the Bayesian MCMC literature) and may, e.g., include microcanonical samplers [67, 68]. However, the requirements of generative diffusion models are more strict and demand that the forward process's asymptotic joint distribution over $\mathbf{x}$ and $\mathbf{m}$ has a simple form that enables sampling in constant time. In contrast, Bayesian MCMC only requires the $\mathbf{x}$-marginal to converge to the prescribed posterior. We still think that relaxing the requirements on tractability enables potentially promising new samplers for future exploration.

# References

[1] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.

[2] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.

[3] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.

[4] 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*, .

[5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.

[6] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*, 23(47):1–33, 2022.

[7] 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*, pages 10684–10695, 2022.

[8] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022. URL `https://arxiv.org/abs/2204.06125`.

[9] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*.

[10] Ruihan Yang, Prakhar Srivastava, and Stephan Mandt. Diffusion probabilistic modeling for video generation, 2022. URL `https://arxiv.org/abs/2203.09481`.

[11] Jonathan Ho, Tim Salimans, Alexey A Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. In *Advances in Neural Information Processing Systems*.

[12] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2837–2845, 2021.

[13] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5826–5835, 2021.

[14] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Score-based generative modeling with critically-damped langevin diffusion. In *International Conference on Learning Representations*, .

[15] Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ICML'11, page 681–688, Madison, WI, USA, 2011. Omnipress. ISBN 9781450306195.

[16] Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient hamiltonian monte carlo. In *International conference on machine learning*, pages 1683–1691. PMLR, 2014.

[17] Yi-An Ma, Tianqi Chen, and Emily Fox. A complete recipe for stochastic gradient mcmc. *Advances in neural information processing systems*, 28, 2015.

[18] Alex Krizhevsky. Learning multiple layers of features from tiny images. pages 32–33, 2009. URL https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf.

[19] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.

[20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.

[21] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016.

[22] Brian D.O. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982. ISSN 0304-4149. doi: https://doi.org/10.1016/0304-4149(82)90051-5. URL https://www.sciencedirect.com/science/article/pii/0304414982900515.

[23] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011. doi: 10.1162/NECO_a_00142.

[24] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. URL https://openreview.net/forum?id=AklttWFnxS9.

[25] Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng, editors. *Handbook of Markov Chain Monte Carlo*. Chapman and Hall/CRC, may 2011. doi: 10.1201/b10905. URL https://doi.org/10.1201%2Fb10905.

[26] L Yin and P Ao. Existence and construction of dynamical potential in nonequilibrium processes without detailed balance. *Journal of Physics A: Mathematical and General*, 39(27):8593, jun 2006. doi: 10.1088/0305-4470/39/27/003. URL https://dx.doi.org/10.1088/0305-4470/39/27/003.

[27] Raghav Singhal, Mark Goldstein, and Rajesh Ranganath. Where to diffuse, how to diffuse, and how to get back: Automated learning for multivariate diffusions. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=osei3IzUia.

[28] Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in Artificial Intelligence*, pages 574–584. PMLR, 2020.

[29] Simo Särkkä and Arno Solin. *Applied stochastic differential equations*, volume 10. Cambridge University Press, 2019.

[30] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems*, 34:11287–11302, 2021.

[31] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *Advances in Neural Information Processing Systems*.

[32] Peter E. Kloeden and Eckhard Platen. *Numerical Solution of Stochastic Differential Equations*. Springer Berlin Heidelberg, 1992. doi: 10.1007/978-3-662-12616-5. URL https://doi.org/10.1007/978-3-662-12616-5.

[33] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197, 2020.

[34] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.

[35] Kushagra Pandey, Avideep Mukherjee, Piyush Rai, and Abhishek Kumar. Diffusevae: Efficient, controllable and high-fidelity generation from low-dimensional latents. *Transactions on Machine Learning Research*.

[36] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.

[37] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.

[38] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=PqvMRDCJT9t`.

[39] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, .

[40] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[41] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. In *International Conference on Learning Representations*.

[42] Eliya Nachmani, Robin San Roman, and Lior Wolf. Non gaussian denoising diffusion models. *arXiv preprint arXiv:2106.07582*, 2021.

[43] J.R. Dormand and P.J. Prince. A family of embedded runge-kutta formulae. *Journal of Computational and Applied Mathematics*, 6(1):19–26, 1980. ISSN 0377-0427. doi: https://doi.org/10.1016/0771-050X(80)90013-3. URL `https://www.sciencedirect.com/science/article/pii/0771050X80900133`.

[44] Bowen Jing, Gabriele Corso, Renato Berlinghieri, and Tommi Jaakkola. Subspace diffusion generative models. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXIII*, pages 274–289. Springer, 2022.

[45] Alexia Jolicoeur-Martineau, Rémi Piché-Taillefer, Ioannis Mitliagkas, and Remi Tachet des Combes. Adversarial score matching and improved sampling for image generation. In *International Conference on Learning Representations*, 2021. URL `https://openreview.net/forum?id=eLfqMl3z3lq`.

[46] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*.

[47] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

[48] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation. In *International Conference on Learning Representations*.

[49] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob Mcgrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, pages 16784–16804. PMLR, 2022.

[50] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

[51] Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*, 479:47–59, 2022.

[52] Ruihan Yang, Prakhar Srivastava, and Stephan Mandt. Diffusion probabilistic modeling for video generation. *arXiv preprint arXiv:2203.09481*, 2022.

[53] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *Transactions on Machine Learning Research*.

[54] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

[55] B. Leimkuhler. *Molecular dynamics : with deterministic and stochastic numerical methods / Ben Leimkuhler, Charles Matthews.* Interdisciplinary applied mathematics, 39. Springer, Cham, 2015. ISBN 3319163744.

[56] Abhishek Sinha, Jiaming Song, Chenlin Meng, and Stefano Ermon. D2c: Diffusion-decoding models for few-shot conditional generation. *Advances in Neural Information Processing Systems*, 34:12533–12548, 2021.

[57] Arash Vahdat and Jan Kautz. Nvae: A deep hierarchical variational autoencoder. *Advances in neural information processing systems*, 33:19667–19679, 2020.

[58] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021.

[59] Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-DPM: an analytic estimate of the optimal reverse variance in diffusion probabilistic models. In *International Conference on Learning Representations*, 2022. URL `https://openreview.net/forum?id=0xiJLKH-ufZ`.

[60] Zhifeng Kong and Wei Ping. On fast sampling of diffusion probabilistic models. In *ICML Workshop on Invertible Neural Networks, Normalizing Flows, and Explicit Likelihood Models*.

[61] Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential integrator. In *The Eleventh International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=Loek7hfb46P`.

[62] Qinsheng Zhang, Molei Tao, and Yongxin Chen. gddim: Generalized denoising diffusion implicit models. *arXiv preprint arXiv:2206.05564*, 2022.

[63] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Genie: Higher-order denoising diffusion solvers. In *Advances in Neural Information Processing Systems*, .

[64] Chenlin Meng, Ruiqi Gao, Diederik P Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *NeurIPS 2022 Workshop on Score-Based Methods*.

[65] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021.

[66] Weitao Du, Tao Yang, He Zhang, and Yuanqi Du. A flexible diffusion model, 2022.

[67] Jakob Robnik and Uroš Seljak. Microcanonical langevin monte carlo. *arXiv preprint arXiv:2303.18221*, 2023.

[68] Greg Ver Steeg and Aram Galstyan. Hamiltonian dynamics with non-newtonian momentum for rapid sampling. *Advances in Neural Information Processing Systems*, 34:11012–11025, 2021.

[69] H. F. Trotter. On the product of semi-groups of operators. *Proceedings of the American Mathematical Society*, 10(4):545–551, 1959. doi: 10.1090/s0002-9939-1959-0108732-6. URL `https://doi.org/10.1090/s0002-9939-1959-0108732-6`.

[70] Gilbert Strang. On the construction and comparison of difference schemes. *SIAM Journal on Numerical Analysis*, 5(3):506–517, September 1968. doi: 10.1137/0705041. URL `https://doi.org/10.1137/0705041`.

[71] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015.

[72] Richard Zhang. Making convolutional networks shift-invariant again. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7324–7334. PMLR, 09–15 Jun 2019. URL `https://proceedings.mlr.press/v97/zhang19a.html`.

[73] Ricky T. Q. Chen. torchdiffeq, 2018. URL `https://github.com/rtqichen/torchdiffeq`.

[74] Anton Obukhov, Maximilian Seitzer, Po-Wei Wu, Semen Zhydenko, Jonathan Kyl, and Elvis Yu-Jing Lin. High-fidelity performance metrics for generative models in pytorch, 2020. URL `https://github.com/toshas/torch-fidelity`. Version: 0.3.0, DOI: 10.5281/zenodo.4957738.

# Contents

# A  A Complete Recipe for SGMs

## A.1  Proof of Theorems

### A.1.1  Proof of Stationarity

Given a positive semi-definite diffusion matrix $\boldsymbol{D}(\mathbf{z})$ and a skew-symmetric matrix $\boldsymbol{Q}(\mathbf{z})$, we can parameterize the drift $f(\mathbf{z})$ for a stochastic process: $d\mathbf{z} = \boldsymbol{f}(\mathbf{z})dt + \sqrt{2\boldsymbol{D}(\mathbf{z})}d\mathbf{w}_t$ as follows:

$$\boldsymbol{f}(\mathbf{z}) = -(\boldsymbol{D}(\mathbf{z}) + \boldsymbol{Q}(\mathbf{z}))\nabla H + \tau(\mathbf{z}), \qquad \tau_i(\mathbf{z}) = \sum_j \frac{\partial}{\partial \mathbf{z}_j}\left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right) \tag{18}$$

Theorem 2.1 then states that the distribution $p_s(z) \propto \exp\left(-H(\mathbf{z})\right)$ will be the stationary distribution for the stochastic process as defined above.

*Proof.* Using the Fokker-Planck formulation for the stochastic dynamics, we have:

$$\frac{\partial p_t(\mathbf{z})}{\partial t} = -\sum_i \frac{\partial}{\partial z_i}\left(\boldsymbol{f}_i(\mathbf{z})p_t(\mathbf{z})\right) + \sum_{i,j} \frac{\partial^2}{\partial \mathbf{z}_i \mathbf{z}_j}\left(\boldsymbol{D}_{ij}(\mathbf{z})p_t(\mathbf{z})\right) \tag{19}$$

Furthermore, we have:

$$\boldsymbol{f}_i(\mathbf{z}) = \tau_i(\mathbf{z}) - \sum_j \left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)\nabla H_j(\mathbf{z}) \tag{20}$$

Therefore,

$$\sum_i \frac{\partial}{\partial \mathbf{z}_i}\left(\boldsymbol{f}_i(\mathbf{z})p_t(\mathbf{z})\right) = \sum_i \frac{\partial}{\partial \mathbf{z}_i}\left[\tau_i(\mathbf{z})p_t(\mathbf{z}) - \sum_j \left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)\nabla H_j(\mathbf{z})p_t(\mathbf{z})\right] \tag{21}$$

$$= \sum_i \frac{\partial}{\partial \mathbf{z}_i}\left[\sum_j \frac{\partial}{\partial \mathbf{z}_j}\left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)p_t(\mathbf{z}) - \sum_j \left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)\nabla H_j(\mathbf{z})p_t(\mathbf{z})\right] \tag{22}$$

$$= \sum_{i,j} \frac{\partial}{\partial \mathbf{z}_i}\left[\frac{\partial}{\partial \mathbf{z}_j}\left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)p_t(\mathbf{z})\right] - \sum_{i,j} \frac{\partial}{\partial \mathbf{z}_i}\left[\left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)\nabla H_j(\mathbf{z})p_t(\mathbf{z})\right] \tag{23}$$

$$= \sum_{i,j} \frac{\partial}{\partial \mathbf{z}_i}\left[\frac{\partial}{\partial \mathbf{z}_j}\left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)p_t(\mathbf{z})\right] - \underbrace{\sum_{i,j} \frac{\partial}{\partial \mathbf{z}_i}\left[\left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)\nabla H_j(\mathbf{z})p_t(\mathbf{z})\right]}_{=F(\mathbf{z})} \tag{24}$$

Substituting the above result in the Fokker-Planck formulation in Eqn. 19, we have:

$$\frac{\partial p_t(\mathbf{z})}{\partial t} = F(\mathbf{z}) - \sum_{i,j} \frac{\partial}{\partial \mathbf{z}_i}\left[\frac{\partial}{\partial \mathbf{z}_j}\left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)p_t(\mathbf{z}) - \frac{\partial}{\partial \mathbf{z}_j}\left(\boldsymbol{D}_{ij}(\mathbf{z})p_t(\mathbf{z})\right)\right] \tag{25}$$

$$= F(\mathbf{z}) - \sum_{i,j} \frac{\partial}{\partial \mathbf{z}_i}\left[\frac{\partial}{\partial \mathbf{z}_j}\left(\boldsymbol{Q}_{ij}(\mathbf{z})\right)p_t(\mathbf{z}) - \boldsymbol{D}_{ij}(\mathbf{z})\frac{\partial}{\partial \mathbf{z}_j}\left(p_t(\mathbf{z})\right)\right] \tag{26}$$

$$= F(\mathbf{z}) - \sum_{i,j} \frac{\partial}{\partial \mathbf{z}_i}\left[\frac{\partial}{\partial \mathbf{z}_j}\left(\boldsymbol{Q}_{ij}(\mathbf{z})p_t(\mathbf{z})\right) - \left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)\frac{\partial}{\partial \mathbf{z}_j}\left(p_t(\mathbf{z})\right)\right] \tag{27}$$

$$= F(\mathbf{z}) + \underbrace{\sum_{i,j} \frac{\partial}{\partial \mathbf{z}_i}\left[\left(\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})\right)\frac{\partial p_t(\mathbf{z})}{\partial \mathbf{z}_j}\right]}_{=G(\mathbf{z})} - \sum_{i,j} \frac{\partial^2}{\partial \mathbf{z}_i \partial \mathbf{z}_j}\left(\boldsymbol{Q}_{ij}(\mathbf{z})p_t(\mathbf{z})\right) \tag{28}$$

$$= F(\mathbf{z}) + G(\mathbf{z}) - \sum_{i,j} \frac{\partial^2}{\partial \mathbf{z}_i \partial \mathbf{z}_j}\left(\boldsymbol{Q}_{ij}(\mathbf{z})p_t(\mathbf{z})\right) \tag{29}$$

Since $\boldsymbol{Q}(\mathbf{z})$ is a skew-symmetric matrix, $\sum_{i,j} \frac{\partial^2}{\partial \mathbf{z}_i \partial \mathbf{z}_j} \left( \boldsymbol{Q}_{ij}(\mathbf{z}) p_t(\mathbf{z}) \right) = 0$. Therefore,

$$\frac{\partial p_t(\mathbf{z})}{\partial t} = F(\mathbf{z}) + G(\mathbf{z}) \tag{30}$$

$$= \sum_{i,j} \frac{\partial}{\partial \mathbf{z}_i} \left[ (\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})) \left( \nabla H_j(\mathbf{z}) p_t(\mathbf{z}) + \frac{\partial p_t(\mathbf{z})}{\partial \mathbf{z}_j} \right) \right] \tag{31}$$

$$= \sum_i \frac{\partial}{\partial \mathbf{z}_i} \left[ (\boldsymbol{D}_i(\mathbf{z}) + \boldsymbol{Q}_i(\mathbf{z})) \left( \nabla H(\mathbf{z}) p_t(\mathbf{z}) + \frac{\partial p_t(\mathbf{z})}{\partial \mathbf{z}} \right) \right] \tag{32}$$

$$= \nabla \cdot \left[ (\boldsymbol{D}_i(\mathbf{z}) + \boldsymbol{Q}_i(\mathbf{z})) \left( \nabla H(\mathbf{z}) p_t(\mathbf{z}) + \frac{\partial p_t(\mathbf{z})}{\partial \mathbf{z}} \right) \right] \tag{33}$$

Therefore, we have the following parameterization for the Fokker-Planck formulation for the defined stochastic dynamics:

$$\frac{\partial p_t(\mathbf{z})}{\partial t} = \nabla \cdot \left[ (\boldsymbol{D}_i(\mathbf{z}) + \boldsymbol{Q}_i(\mathbf{z})) \left( \nabla H(\mathbf{z}) p_t(\mathbf{z}) + \frac{\partial p_t(\mathbf{z})}{\partial \mathbf{z}} \right) \right] \tag{34}$$

Substituting $p_s(\mathbf{z}) \propto \exp(-H(\mathbf{z}))$ in the above result implies $\frac{\partial p_t(\mathbf{z})}{\partial t} = 0$. This implies that $p_s(\mathbf{z}) \propto \exp(-H(\mathbf{z}))$ is the form of the stationary distribution for the drift parameterization $\boldsymbol{f}(\mathbf{z}) = -(\boldsymbol{D}(\mathbf{z}) + \boldsymbol{Q}(\mathbf{z}))\nabla H + \boldsymbol{\tau}(\mathbf{z})$. An alternative version of this proof can be found in Ma et al. [17] □

### A.1.2 Proof of Completeness

We now state the proof for Theorem 2.2 which states that for every stochastic dynamics $d\mathbf{z} = \boldsymbol{f}(\mathbf{z})dt + \sqrt{2\boldsymbol{D}(\mathbf{z})}d\mathbf{w}_t$ with the desired stationary distribution $p_s(\mathbf{z}) \propto \exp(-H(\mathbf{z}))$, there exists a positive semi-definite $\boldsymbol{D}(\mathbf{z})$ and a skew-symmetric $\boldsymbol{Q}(\mathbf{z})$ such that $\boldsymbol{f}(\mathbf{z}) = -(\boldsymbol{D}(\mathbf{z}) + \boldsymbol{Q}(\mathbf{z}))\nabla H + \boldsymbol{\tau}(\mathbf{z})$ holds. We directly include the proof from Ma et al. [17] for completeness.

*Proof.* We have the following result:

$$\boldsymbol{f}_i(\mathbf{z})p_s(\mathbf{z}) = \boldsymbol{\tau}_i(\mathbf{z})p_s(\mathbf{z}) - \sum_j (\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})) \nabla H_j(\mathbf{z})p_s(\mathbf{z}) \tag{35}$$

$$= \sum_j \frac{\partial}{\partial \mathbf{z}_j} (\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})) p_s(\mathbf{z}) - \sum_j (\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})) \nabla H_j(\mathbf{z})p_s(\mathbf{z}) \tag{36}$$

$$= \sum_j \frac{\partial}{\partial \mathbf{z}_j} \left[ (\boldsymbol{D}_{ij}(\mathbf{z}) + \boldsymbol{Q}_{ij}(\mathbf{z})) p_s(\mathbf{z}) \right] \tag{37}$$

which implies,

$$\sum_j \frac{\partial}{\partial \mathbf{z}_j} (\boldsymbol{Q}_{ij}(\mathbf{z})p_s(\mathbf{z})) = \boldsymbol{f}_i(\mathbf{z})p_s(\mathbf{z}) - \sum_j \frac{\partial}{\partial \mathbf{z}_j} (\boldsymbol{D}_{ij}(\mathbf{z})p_s(\mathbf{z})) \tag{38}$$

Furthermore, from the Fokker-Planck formalism,

$$\frac{\partial p_t(\mathbf{z})}{\partial t} = -\sum_i \frac{\partial}{\partial z_i} (\boldsymbol{f}_i(\mathbf{z})p_t(\mathbf{z})) + \sum_{i,j} \frac{\partial^2}{\partial \mathbf{z}_i \mathbf{z}_j} (\boldsymbol{D}_{ij}(\mathbf{z})p_t(\mathbf{z})) \tag{39}$$

$$= -\sum_i \frac{\partial}{\partial z_i} \left[ \boldsymbol{f}_i(\mathbf{z})p_t(\mathbf{z}) - \sum_j \frac{\partial}{\partial \mathbf{z}_j} (\boldsymbol{D}_{ij}(\mathbf{z})p_t(\mathbf{z})) \right] \tag{40}$$

For $p_t(\mathbf{z}) = p_s(\mathbf{z})$, we have,

$$\sum_i \frac{\partial}{\partial z_i} \left[ \boldsymbol{f}_i(\mathbf{z})p_t(\mathbf{z}) - \sum_j \frac{\partial}{\partial \mathbf{z}_j} (\boldsymbol{D}_{ij}(\mathbf{z})p_t(\mathbf{z})) \right] = 0 \tag{41}$$

20

Denoting the Fourier transform of $\boldsymbol{Q}(\mathbf{z})p_s(\mathbf{z})$ as $\hat{\boldsymbol{Q}}(\boldsymbol{k})$ and the Fourier transform of $\boldsymbol{f}_i(\mathbf{z})p_t(\mathbf{z}) - \sum_j \frac{\partial}{\partial \mathbf{z}_j}(\boldsymbol{D}_{ij}(\mathbf{z})p_t(\mathbf{z}))$ by $\hat{\boldsymbol{F}}(\boldsymbol{k})$, then from Eqns. 38 and 41 we have the following equations in the Fourier-space:

$$2\pi i \hat{\boldsymbol{Q}} \boldsymbol{k} = \hat{\boldsymbol{F}} \tag{42}$$

$$\boldsymbol{k}^T \hat{\boldsymbol{F}} = 0 \tag{43}$$

Therefore, it implies that the matrix $\hat{\boldsymbol{Q}}$ is a projection matrix from $\boldsymbol{k}$ to the span of $\hat{\boldsymbol{F}}$. Consequently, the matrix $\hat{\boldsymbol{Q}}$ can be constructed as: $\hat{\boldsymbol{Q}} = (2\pi i)^{-1} \frac{\hat{\boldsymbol{F}} \boldsymbol{k}^T}{\boldsymbol{k}^T \boldsymbol{k}} - (2\pi i)^{-1} \frac{\boldsymbol{k} \hat{\boldsymbol{F}}^T}{\boldsymbol{k}^T \boldsymbol{k}}$. This construction also shows that $\hat{\boldsymbol{Q}}$ is skew-symmetric. Moreover, the skew-symmetric $\boldsymbol{Q}$ can be obtained by computing the Inverse-Fourier transform of $(p_s(\mathbf{z}))^{-1} \hat{\boldsymbol{Q}}$ $\qquad\square$

## A.2 Existing SGMs parameterized using the SGM recipe

In this section, we provide examples of SGMs that can be cast under the recipe proposed in Section 2.2. It is worth noting that under the completeness framework proposed in Section 2.2, given a positive semi-definite diffusion matrix $\boldsymbol{D}(\mathbf{z})$, a skew-symmetric curl matrix $\boldsymbol{Q}(\mathbf{z})$ and the Hamiltonian $H(\mathbf{z})$ corresponding to a specified target distribution $p_s(\mathbf{z})$, the forward process SDE can be parameterized in terms of the target distribution as follows:

$$H(\mathbf{z}) = U(\mathbf{x}) + \frac{\mathbf{m}^T M^{-1} \mathbf{m}}{2}, \quad \nabla H(\mathbf{z}) = \begin{pmatrix} \nabla U(\mathbf{x}) \\ M^{-1}\mathbf{m} \end{pmatrix} \tag{44}$$

$$\boldsymbol{f}(\mathbf{z}) = -(\boldsymbol{D}(\mathbf{z}) + \boldsymbol{Q}(\mathbf{z})) \begin{pmatrix} \nabla U(\mathbf{x}) \\ M^{-1}\mathbf{m} \end{pmatrix} + \tau(\mathbf{z}) \tag{45}$$

$$d\mathbf{z} = \boldsymbol{f}(\mathbf{z})dt + \sqrt{2\boldsymbol{D}(\mathbf{z})}d\mathbf{w}_t \tag{46}$$

We now recast several existing SGMs under this framework:

### A.2.1 Non-augmented SGMs

For SGMs with a non-augmented form, we assume auxiliary variables $\mathbf{m}_t = 0$ with the equilibrium distribution given by $p_s(\mathbf{z}) = \mathcal{N}(\mathbf{0_d}, \mathbf{I_d})$. The Hamiltonian and its gradient can then be specified as follows:

$$H(\mathbf{z}) = \frac{\mathbf{x}^T \mathbf{x}}{2}, \qquad \nabla H(\mathbf{z}) = \mathbf{x} \tag{47}$$

For the choice of $\boldsymbol{D}_{\text{VP}}(\mathbf{z}) = \frac{\beta_t}{2}\boldsymbol{I}_d$ and $\boldsymbol{Q}_{\text{VP}}(\mathbf{z}) = \mathbf{0}_d$, the drift for the forward SDE defined in Eqn. 45 reduces to the following form:

$$d\mathbf{x} = -\frac{\beta_t}{2}\mathbf{x}dt + \sqrt{\beta_t}d\mathbf{w}_t \tag{48}$$

where $\beta_t$ is a time-dependent constant. The forward SDE in Eqn. 48 is the same as the VP-SDE proposed in [4]. From our recipe, the stationary distribution for the VPSDE should be $p_s(\mathbf{x}) = \mathcal{N}(\mathbf{0}_d, \boldsymbol{I}_d)$. Indeed the perturbation kernel for the VP-SDE is specified as follows:

$$p(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_0 e^{-\frac{1}{2}\int_0^t \beta(s)ds}, (1 - e^{-\int_0^t \beta(s)ds})^2 \boldsymbol{I}_d) \tag{49}$$

which converges to a standard Gaussian distribution as $t \to \infty$. This example suggests that the proposed recipe can be used to establish the validity of the convergence of a forward process with a specified stationary distribution $p_s(\mathbf{z})$ without deriving the perturbation kernel or relying on physical intuition. Interestingly, the Variance-Exploding (VE) SDE [4] is one example that cannot be cast in our framework. This would mean that it will not asymptotically converge to the standard Gaussian distribution at equilibrium. Indeed, this can be confirmed from the analytical form of the perturbation kernel of the VE-SDE given by:

$$p(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_0, [\sigma^2(t) - \sigma^2(0)]\boldsymbol{I}_d) \tag{50}$$

As $t \to \infty$, the variance of the perturbation kernel of the VE-SDE grows unbounded and therefore does not converge to the equilibrium distribution $\mathcal{N}(\mathbf{0}_d, \boldsymbol{I}_d)$. This should not be surprising since the VE-SDE, for the specified Hamiltonian, could not be recast in the completeness framework, to begin with.

### A.2.2 Augmented SGMs

For SGMs with an augmented state-space (data state space $\mathbf{x}_t$ + auxiliary variables $\mathbf{m}_t$), we assume the equilibrium distribution $p_s(\mathbf{z}) = \mathcal{N}(\mathbf{0_d}, \mathbf{I_d})\mathcal{N}(\mathbf{0_d}, \mathbf{MI_d})$. The Hamiltonian and its gradient can then be specified as follows:

$$H(\mathbf{z}) = \frac{\mathbf{x}^T\mathbf{x}}{2} + \frac{\mathbf{m}^T M^{-1}\mathbf{m}}{2}, \qquad \nabla H(\mathbf{z}) = \begin{pmatrix} \mathbf{x} \\ M^{-1}\mathbf{m} \end{pmatrix} \tag{51}$$

For this choice of $H$, the forward SDE representative of PSLD can be obtained by choosing the following $\boldsymbol{D}$ and $\boldsymbol{Q}$ matrices:

$$\boldsymbol{D}_{\text{PSLD}} = \frac{\beta}{2}\left(\begin{pmatrix} \Gamma & 0 \\ 0 & M\gamma \end{pmatrix} \otimes \boldsymbol{I}_d\right) \qquad \boldsymbol{Q}_{\text{PSLD}} = \frac{\beta}{2}\left(\begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \otimes \boldsymbol{I}_d\right) \tag{52}$$

Similarly, the forward SDE representative of CLD [14] can be obtained by choosing:

$$\boldsymbol{D}_{\text{CLD}} = \beta\left(\begin{pmatrix} 0 & 0 \\ 0 & \Gamma \end{pmatrix} \otimes \boldsymbol{I}_d\right) \qquad \boldsymbol{Q}_{\text{CLD}} = \beta\left(\begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \otimes \boldsymbol{I}_d\right) \tag{53}$$

Since both PSLD and CLD can be shown to converge asymptotically to $p_s(\mathbf{z})$ from the analytical form of their perturbation kernels $p(\mathbf{z}_t|\mathbf{z}_0)$, the result from our completeness framework is valid. More importantly, given a forward process for an SGM, our recipe can be used to validate if the SGM converges to a specified equilibrium distribution without the need for analytically determining the perturbation kernel (which is usually non-trivial).

## B  Phase Space Langevin Diffusion

In this section, we elaborate on several aspects of PSLD, which were discussed briefly in the main text. Moreover, we work with the following form of the forward process for PSLD:

$$\begin{pmatrix} d\mathbf{x}_t \\ d\mathbf{m}_t \end{pmatrix} = \left(\frac{\beta_t}{2}\begin{pmatrix} -\Gamma & M^{-1} \\ -1 & -\nu \end{pmatrix} \otimes \boldsymbol{I}_d\right)\begin{pmatrix} \mathbf{x}_t \\ \mathbf{m}_t \end{pmatrix} dt + \left(\begin{pmatrix} \sqrt{\Gamma\beta_t} & 0 \\ 0 & \sqrt{M\nu\beta_t} \end{pmatrix} \otimes \boldsymbol{I}_d\right) d\mathbf{w}_t, \tag{54}$$

It is worth noting that the form of the forward process defined in Eqn. 54 is more general than Eqn. 12 in the sense that we consider a time-dependent $\beta_t$ here for our discussions. We can then reason about the forward SDE in Eqn. 12 by fixing $\beta_t$ to a time-independent quantity $\beta$ for all subsequent analyses.

### B.1  Critical Damping in PSLD

Assuming $\beta_t = 1$ and $\mathbf{x}_t, \mathbf{m}_t \in \mathbb{R}$ for simplicity, the equations of motion for the deterministic dynamics can be specified as follows:

$$\frac{dx_t}{dt} = -\Gamma x_t + M^{-1}m_t \tag{55}$$

$$\frac{dm_t}{dt} = -x_t - \nu m_t \tag{56}$$

From Eqn. 55, we have:

$$m_t = M\left(\frac{dx_t}{dt} + \Gamma x_t\right) \tag{57}$$

Furthermore, taking the derivative of both sides in Eqn. 55, we have:

$$\frac{d^2x_t}{dt^2} = -\Gamma\frac{dx_t}{dt} + M^{-1}\frac{dm_t}{dt} \tag{58}$$

$$= -\Gamma\frac{dx_t}{dt} + M^{-1}[-x_t - \nu m_t] \tag{59}$$

$$= -\Gamma\frac{dx_t}{dt} + M^{-1}\left[-x_t - \nu M\left(\frac{dx_t}{dt} + \Gamma x_t\right)\right] \tag{60}$$

$$= -\Gamma \frac{dx_t}{dt} - M^{-1}x_t - \nu \left( \frac{dx_t}{dt} + \Gamma x_t \right) \tag{61}$$

$$= -\Gamma \frac{dx_t}{dt} - M^{-1}x_t - \nu \frac{dx_t}{dt} - \Gamma \nu x_t \tag{62}$$

$$= -(\Gamma + \nu) \frac{dx_t}{dt} - M^{-1}x_t - \Gamma \nu x_t \tag{63}$$

We, therefore, have the following dynamical equation in terms of the position:

$$\frac{d^2 x_t}{dt^2} + (\Gamma + \nu) \frac{dx_t}{dt} + (M^{-1} + \Gamma \nu)x_t = 0 \tag{64}$$

Assuming the exponential *ansatz* $\mathbf{x}_t = \exp(-\lambda t)$ and plugging into the above ODE, we have the following result:

$$\exp(-\lambda t) \left[ \lambda^2 - (\Gamma + \nu)\lambda + (M^{-1} + \Gamma \nu) \right] = 0 \tag{65}$$

which implies,

$$\lambda^2 - (\Gamma + \nu)\lambda + (M^{-1} + \Gamma \nu) = 0 \tag{66}$$

$$\lambda = \frac{(\Gamma + \nu) \pm \sqrt{(\Gamma + \nu)^2 - 4M^{-1} - 4\Gamma \nu}}{2} \tag{67}$$

$$\lambda = \frac{(\Gamma + \nu) \pm \sqrt{(\Gamma - \nu)^2 - 4M^{-1}}}{2} \tag{68}$$

Corresponding to the value of $\nu, \Gamma$ and $M$, we can now have the following damping conditions:

**(i)** $(\Gamma - \nu)^2 < 4M^{-1}$ corresponds to *Underdamped dynamics*

**(ii)** $(\Gamma - \nu)^2 = 4M^{-1}$ corresponds to *Critical damping*

**(iii)** $(\Gamma - \nu)^2 > 4M^{-1}$ corresponds to *Overdamped dynamics*

Moreover when $\Gamma = 0$ and $\bar{\nu} = M\nu$, we get: $\bar{\nu}^2 = 4M$ which is the critical damping condition proposed in Dockhorn et al. [14]. Therefore, similar to Dockhorn et al. [14], we work in the Critical Damping regime specified by the condition $(\Gamma - \nu)^2 = 4M^{-1}$

### B.2 PSLD Training

#### B.2.1 Overall Training Framework in PSLD

Following the derivation in Dockhorn et al. [14], the maximum likelihood training formulation for score matching can be specified as follows. Let $p_0, q_0$ be two densities with corresponding marginal densities $p_t$ and $q_t$ (for forward diffusion using PSLD defined in Eqn. 54) at time t. As shown in Song et al. [24], the KL-Divergence between $p_0$ and $q_0$ can then be expressed as a mixture of score-matching losses over multiple time scales as follows:

$$D_{\mathrm{KL}}(p_0 \parallel q_0) = D_{\mathrm{KL}}(p_0 \parallel q_0) - D_{\mathrm{KL}}(p_T \parallel q_T) + D_{\mathrm{KL}}(p_T \parallel q_T)$$
$$= -\int_0^T \frac{\partial D_{\mathrm{KL}}(p_t \parallel q_t)}{\partial t} dt + D_{\mathrm{KL}}(p_T \parallel q_T) \tag{69}$$

Following the derivation from Song et al. [24], the Fokker-Planck equation describing the time evolution of the probability density function of the SDE in Eqn. 54 can be expressed as follows:

$$\frac{\partial p_t(\mathbf{z}_t)}{\partial t} = \nabla_{\mathbf{z}_t} \cdot \left[ \frac{1}{2} \left( G(t)G(t)^\top \otimes \mathbf{I}_d \right) \nabla_{\mathbf{z}_t} p_t(\mathbf{z}_t) - p_t(\mathbf{z}_t)(f(t) \otimes \mathbf{I}_d)\mathbf{z}_t \right]$$
$$= \nabla_{\mathbf{z}_t} \cdot \left[ \mathbf{h}_p(\mathbf{z}_t, t) p_t(\mathbf{z}_t) \right] \tag{70}$$

where

$$\mathbf{h}_p(\mathbf{z}_t, t) := \frac{1}{2} \left( G(t)G(t)^\top \otimes \mathbf{I}_d \right) \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) - (f(t) \otimes \mathbf{I}_d)\mathbf{z}_t \tag{71}$$

$$f(t) = \left( \frac{\beta_t}{2} \begin{pmatrix} -\Gamma & M^{-1} \\ -1 & -\nu \end{pmatrix} \right) \tag{72}$$

23

$$G(t) = \begin{pmatrix} \sqrt{\Gamma \beta_t} & 0 \\ 0 & \sqrt{M\nu\beta_t} \end{pmatrix} \tag{73}$$

Further assuming that $\log p_t(\mathbf{z}_t)$ and $\log q_t(\mathbf{z}_t)$ are smooth functions with at most polynomial growth at infinity, we have

$$\lim_{\mathbf{z}_t \to \infty} \boldsymbol{h}_p(\mathbf{z}_t, t) p_t(\mathbf{z}_t) = \lim_{\mathbf{z}_t \to \infty} \boldsymbol{h}_q(\mathbf{z}_t, t) q_t(\mathbf{z}_t) = 0. \tag{74}$$

Using the above fact, we can compute the time-derivative of the Kullback–Leibler divergence between $p_t$ and $q_t$ as

$$
\begin{aligned}
\frac{\partial D_{\mathrm{KL}}(p_t \parallel q_t)}{\partial t} &= \frac{\partial}{\partial t} \int p_t(\mathbf{z}_t) \log \frac{p_t(\mathbf{z}_t)}{q_t(\mathbf{z}_t)} \, d\mathbf{z}_t \\
&= \int \frac{\partial p_t(\mathbf{z}_t)}{\partial t} \log \frac{p_t(\mathbf{z}_t)}{q_t(\mathbf{z}_t)} \, d\mathbf{z}_t - \int \frac{p_t(\mathbf{z}_t)}{q_t(\mathbf{z}_t)} \frac{\partial q_t(\mathbf{z}_t)}{\partial t} d\mathbf{z}_t \\
&= -\int p_t(\mathbf{z}_t) \Big[ \boldsymbol{h}_p(\mathbf{z}_t, t) - \boldsymbol{h}_q(\mathbf{z}_t, t) \Big]^\top \Big[ \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{z}_t} \log q_t(\mathbf{z}_t) \Big] d\mathbf{z}_t \\
&= -\frac{1}{2} \int p_t(\mathbf{z}_t) \Big[ \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{z}_t} \log q_t(\mathbf{z}_t) \Big]^\top \Big( G(t)G(t)^\top \otimes \boldsymbol{I}_d \Big) \\
&\qquad\qquad \Big[ \nabla_{\mathbf{u}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{z}_t} \log q_t(\mathbf{z}_t) \Big] d\mathbf{z}_t \\
&= -\frac{1}{2} \int p_t(\mathbf{z}_t) \Big[ \Gamma\beta_t \| \nabla_{\mathbf{x}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{z}_t) \|_2^2 + \\
&\qquad\qquad M\nu\beta_t \| \nabla_{\mathbf{m}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{m}_t} \log q_t(\mathbf{z}_t) \|_2^2 \Big] d\mathbf{z}_t
\end{aligned}
\tag{75}
$$

Assuming our generative prior $p(x_T)$ matches the equilibrium state of the forward process closely i.e. $D_{\mathrm{KL}}(p_T \parallel q_T) \approx 0$ and substituting the result in Eqn. 75 in Eqn. 69, we get the following score-matching objective corresponding to the maximum-likelihood objective in Eqn. 69 as follows:

$$
D_{\mathrm{KL}}(p_0 \parallel q_0) = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{z}_t \sim p_t(\mathbf{z}_t)} \Big[ \underbrace{\Gamma\beta_t \big\| \nabla_{\mathbf{x}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{z}_t) \big\|_2^2}_{\text{Data-Space}} + \\
\underbrace{M\nu\beta_t \big\| \nabla_{\mathbf{m}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{m}_t} \log q_t(\mathbf{z}_t) \big\|_2^2}_{\text{Momentum-Space}} \Big]
\tag{76}
$$

In general, the above score-matching loss can be re-formulated using arbitrary loss weightings $\lambda_1(t)$ and $\lambda_2(t)$ as follows:

$$
D_{\mathrm{KL}}(p_0 \parallel q_0) = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{z}_t \sim p_t(\mathbf{z}_t)} \Big[ \lambda_1(t) \big\| \nabla_{\mathbf{x}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{z}_t) \big\|_2^2 + \\
\lambda_2(t) \big\| \nabla_{\mathbf{m}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{m}_t} \log q_t(\mathbf{z}_t) \big\|_2^2 \Big]
\tag{77}
$$

Choosing the same weighting for both loss components i.e. $\lambda_1(t) = \lambda_2(t) = \lambda(t)$, the score-matching objective in Eqn. 76 can be simplified as follows:

$$
\mathcal{L}_{\mathrm{SM}} = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{z}_t \sim p_t(\mathbf{z}_t)} \Big[ \lambda(t) \big\| \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) - \nabla_{\mathbf{z}_t} \log q_t(\mathbf{z}_t) \big\|_2^2 \Big] \tag{78}
$$

Approximating the score $\nabla_{\mathbf{z}_t} \log q_t(\mathbf{z}_t)$ using a parametric estimator $\boldsymbol{s}_\theta(\mathbf{z}_t, t)$ and following Vincent [23], it can be shown that the $\mathcal{L}_{\mathrm{SM}}$ objective is equivalent to the following Denoising Score Matching (DSM) objective:

$$
\mathcal{L}_{\mathrm{DSM}} = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{z}_0 \sim p(\mathbf{z}_0)} \mathbb{E}_{\mathbf{z}_t \sim p_t(\mathbf{z}_t|\mathbf{z}_0)} \Big[ \lambda(t) \big\| \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t|\mathbf{z}_0) - \boldsymbol{s}_\theta(\mathbf{z}_t, t) \big\|_2^2 \Big] \tag{79}
$$

Moreover, Dockhorn et al. [14] propose to use the following objective a.k.a. Hybrid Score Matching (HSM), which is equivalent to the DSM objective (upto a constant independent of $\theta$):

$$
\mathcal{L}_{\mathrm{HSM}} = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0)} \mathbb{E}_{\mathbf{z}_t \sim p_t(\mathbf{z}_t|\mathbf{x}_0)} \Big[ \lambda(t) \big\| \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t|\mathbf{x}_0) - \boldsymbol{s}_\theta(\mathbf{z}_t, t) \big\|_2^2 \Big] \tag{80}
$$

The perturbation kernels $p(\mathbf{z}_t|\mathbf{z}_0)$ and $p(\mathbf{z}_t|\mathbf{x}_0)$ can be computed analytically for an SDE with affine drift (See Appendix B.3 for the exact analytical forms of the perturbation kernel for PSLD). Following Dockhorn et al. [14], we use the Hybrid Score Matching (HSM) objective throughout this work. We next discuss the computation of the analytical form of the target score $\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t|\mathbf{x}_0)$ (or $\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t|\mathbf{z}_0)$ for DSM) and the parameterization of our score network $\boldsymbol{s}_\theta(\mathbf{z}_t, t)$.

### B.2.2 Analytical Score Computation and Parameterization

In cases when the perturbation kernels are multivariate Gaussian distributions of the form $\mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$, the target score $\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t|\mathbf{x}_0)$ can be computed analytically as follows:

$$\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{x}_0) = -\boldsymbol{\Sigma}_t^{-1}(\mathbf{z}_t - \boldsymbol{\mu}_t) \tag{81}$$

$$= -\boldsymbol{L}_t^{-T}\boldsymbol{L}_t^{-1}(\boldsymbol{L}_t\boldsymbol{\epsilon}) = -\boldsymbol{L}_t^{-T}\boldsymbol{\epsilon} \tag{82}$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}_{2d}, \boldsymbol{I}_{2d})$ and $\boldsymbol{\Sigma}_t = \boldsymbol{L}_t\boldsymbol{L}_t^T$ is the Cholesky decomposition. Moreover, for $\boldsymbol{\Sigma}_t = \left( \begin{pmatrix} \Sigma_t^{xx} & \Sigma_t^{xm} \\ \Sigma_t^{xm} & \Sigma_t^{mm} \end{pmatrix} \otimes \boldsymbol{I}_d \right)$ (as is the case in PSLD), the Cholesky decomposition can be computed analytically as follows:

$$\boldsymbol{L}_t = \left( \begin{pmatrix} L_t^{xx} & 0 \\ L_t^{xm} & L_t^{mm} \end{pmatrix} \otimes \boldsymbol{I}_d \right) \tag{83}$$

$$L_t = \begin{pmatrix} L_t^{xx} & 0 \\ L_t^{xm} & L_t^{mm} \end{pmatrix} = \begin{pmatrix} \sqrt{\Sigma_t^{xx}} & 0 \\ \frac{\Sigma_t^{xm}}{\sqrt{\Sigma_t^{xx}}} & \sqrt{\frac{\Sigma_t^{xx}\Sigma_t^{mm}-(\Sigma_t^{xm})^2}{\Sigma_t^{xx}}} \end{pmatrix} \tag{84}$$

Consequently,

$$\boldsymbol{L}_t^{-T} = L_t^{-T} \otimes \boldsymbol{I}_d$$

$$= \begin{pmatrix} \sqrt{\Sigma_t^{xx}} & \frac{\Sigma_t^{xm}}{\sqrt{\Sigma_t^{xx}}} \\ 0 & \sqrt{\frac{\Sigma_t^{xx}\Sigma_t^{mm}-(\Sigma_t^{xm})^2}{\Sigma_t^{xx}}} \end{pmatrix}^{-1} \otimes \boldsymbol{I}_d$$

$$= \begin{pmatrix} \frac{1}{\sqrt{\Sigma_t^{xx}}} & \frac{-\Sigma_t^{xm}}{\sqrt{\Sigma_t^{xx}}\sqrt{\Sigma_t^{xx}\Sigma_t^{mm}-(\Sigma_t^{xm})^2}} \\ 0 & \frac{\sqrt{\Sigma_t^{xx}\Sigma_t^{mm}-(\Sigma_t^{xm})^2}}{\sqrt{\Sigma_t^{xx}\Sigma_t^{mm}-(\Sigma_t^{xm})^2}} \end{pmatrix} \otimes \boldsymbol{I}_d. \tag{85}$$

Plugging the analytical form of $\boldsymbol{L}_t^{-T}$ into Eqn. 82, we get the following analytical form of the target score $\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t|\mathbf{x}_0)$:

$$\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t|\mathbf{x}_0) = -\boldsymbol{L}_t^{-T}\boldsymbol{\epsilon} \tag{86}$$

$$= -\left( \begin{pmatrix} l_t^{xx} & l_t^{xm} \\ 0 & l_t^{mm} \end{pmatrix} \otimes \boldsymbol{I}_d \right) \begin{pmatrix} \boldsymbol{\epsilon}_x \\ \boldsymbol{\epsilon}_m \end{pmatrix} \tag{87}$$

$$= -\begin{pmatrix} l_t^{xx}\boldsymbol{\epsilon_x} + l_t^{xm}\boldsymbol{\epsilon_m} \\ l_t^{mm}\boldsymbol{\epsilon_m} \end{pmatrix} \tag{88}$$

where $l_t^{xx} = \frac{1}{\sqrt{\Sigma_t^{xx}}}$, $l_t^{xm} = \frac{-\Sigma_t^{xm}}{\sqrt{\Sigma_t^{xx}}\sqrt{\Sigma_t^{xx}\Sigma_t^{mm}-(\Sigma_t^{xm})^2}}$ and $l_t^{mm} = \sqrt{\frac{\Sigma_t^{xx}}{\Sigma_t^{xx}\Sigma_t^{mm}-(\Sigma_t^{xm})^2}}$. While one can directly model the score as defined in Eqn. 88, we instead parameterize the score network as $\boldsymbol{s}_\theta(\mathbf{z}_t, t) = -\boldsymbol{L}_t^{-T}\boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t)$.

### B.2.3 Putting it all together

Plugging the analytical form of $\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t|\mathbf{x}_0)$ and our score network parameterization $\boldsymbol{s}_\theta(\mathbf{z}_t, t) = -\boldsymbol{L}_t^{-T}\boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t)$ in the HSM objective in Eqn. 80, we get the following objective:

$$\mathcal{L}_{\text{HSM}} = \frac{1}{2}\mathbb{E}_{t\sim\mathcal{U}(0,T)}\mathbb{E}_{\mathbf{x}_0\sim p(\mathbf{x}_0)}\mathbb{E}_{\mathbf{z}_t\sim p_t(\mathbf{z}_t|\mathbf{x}_0)}\left[\lambda(t)\big\|\nabla_{\mathbf{z}_t}\log p_t(\mathbf{z}_t|\mathbf{x}_0) - \boldsymbol{s}_\theta(\mathbf{z}_t, t)\big\|_2^2\right] \tag{89}$$

$$= \frac{1}{2}\mathbb{E}_{t\sim\mathcal{U}(0,T)}\mathbb{E}_{\mathbf{x}_0\sim p(\mathbf{x}_0)}\mathbb{E}_{\boldsymbol{\epsilon}\sim\mathcal{N}(\mathbf{0}_{2d},\boldsymbol{I}_{2d})}\left[\lambda(t)\big\|\boldsymbol{L}_t^{-T}\boldsymbol{\epsilon} - \boldsymbol{L}_t^{-T}\boldsymbol{\epsilon}_\theta(\boldsymbol{\mu}_t + \boldsymbol{L}_t\boldsymbol{\epsilon}, t)\big\|_2^2\right] \tag{90}$$

$$\leq \frac{1}{2}\mathbb{E}_{t\sim\mathcal{U}(0,T)}\mathbb{E}_{\mathbf{x}_0\sim p(\mathbf{x}_0)}\mathbb{E}_{\boldsymbol{\epsilon}\sim\mathcal{N}(\mathbf{0}_{2d},\boldsymbol{I}_{2d})}\left[\lambda(t)\big\|\boldsymbol{L}_t^{-T}\big\|_2^2\big\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\boldsymbol{\mu}_t + \boldsymbol{L}_t\boldsymbol{\epsilon}, t)\big\|_2^2\right] \tag{91}$$

It is worth noting that since our original HSM objective is upper bounded by the objective in Eqn. 91, minimizing the latter also minimizes $\mathcal{L}_{\text{HSM}}$. Since we optimize for sample quality, following prior work [4, 14], we choose $\lambda(t) = \frac{1}{\|\boldsymbol{L}_t^{-T}\|_2^2}$ to cancel the weighting induced by $\|\boldsymbol{L}_t^{-T}\|_2^2$. Our final training objective reduces to the following *noise-prediction* formulation:

$$\mathcal{L}(\theta) = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0)} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}_{2d}, \boldsymbol{I}_{2d})} \left[ \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\boldsymbol{\mu}_t + \boldsymbol{L}_t \boldsymbol{\epsilon}, t) \right\|_2^2 \right] \tag{92}$$

It is worth noting that in our training setup, we need to predict the full $2d$-dimensional $\boldsymbol{\epsilon}$. This is in contrast to the training setup in CLD, where we only need to predict the last-d components i.e. $\boldsymbol{\epsilon}_{d:2d}$ of the noise vector $\boldsymbol{\epsilon}$. This difference in training arises due to different formulations of the diffusion coefficient in PSLD and CLD. Indeed, setting $\Gamma = 0$ in Eqn. 76 would result in a similar training objective as in CLD.

### B.3 Perturbation Kernel in PSLD

We now present the analytical form of the perturbation kernels $p(\mathbf{z}_t|\mathbf{z}_0)$ and $p(\mathbf{z}_t|\mathbf{x}_0)$ for PSLD, which are required for training using DSM or HSM respectively.

### B.3.1 Mean and Variance of $p(\mathbf{z}_t|\mathbf{z}_0)$

Since the drift and the diffusion coefficients in Eqn. 54 are affine, the perturbation kernel $p(\mathbf{z}_t|\mathbf{z}_0)$ will be a multivariate Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$. Following Särkkä and Solin [29], $\boldsymbol{\mu}_t$ and $\boldsymbol{\Sigma}_t$, evolve as the following ODEs:

$$\frac{d\boldsymbol{\mu}_t}{dt} = \boldsymbol{F}(t)\boldsymbol{\mu}_t \tag{93}$$

$$\frac{d\boldsymbol{\Sigma}_t}{dt} = \boldsymbol{F}(t)\boldsymbol{\Sigma}_t + \boldsymbol{\Sigma}_t \boldsymbol{F}^T(t) + \boldsymbol{G}(t)\boldsymbol{G}(t)^T \tag{94}$$

where $\boldsymbol{F}(t) = \left( \frac{\beta_t}{2} \begin{pmatrix} -\Gamma & M^{-1} \\ -1 & -\nu \end{pmatrix} \otimes \boldsymbol{I}_d \right)$ and $\boldsymbol{G}(t) = \left( \begin{pmatrix} \sqrt{\Gamma\beta_t} & 0 \\ 0 & \sqrt{M\nu\beta_t} \end{pmatrix} \otimes \boldsymbol{I}_d \right)$ for PSLD. Under critical damping i.e. $M^{-1} = \frac{(\Gamma-\nu)^2}{4}$, solving the ODEs for the mean and variance yields the following form of $p(\mathbf{z}_t|\mathbf{z}_0) = \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$:

$$\boldsymbol{\mu}_t = \begin{pmatrix} \boldsymbol{\mu}_t^x \\ \boldsymbol{\mu}_t^m \end{pmatrix} = \begin{pmatrix} A_1 \mathcal{B}(t)\boldsymbol{x}_0 + A_2 \mathcal{B}(t)\boldsymbol{m}_0 + \boldsymbol{x}_0 \\ C_1 \mathcal{B}(t)\boldsymbol{x}_0 + C_2 \mathcal{B}(t)\boldsymbol{m}_0 + \boldsymbol{m}_0 \end{pmatrix} e^{-\left(\frac{\nu+\Gamma}{4}\right)\mathcal{B}(t)} \tag{95}$$

where $\mathcal{B}(t) = \int_0^t \beta(s)ds$ and coefficients:

$$A_1 = \frac{\nu - \Gamma}{4} \qquad A_2 = \frac{(\Gamma - \nu)^2}{8} \tag{96}$$

$$C_1 = \frac{-1}{2} \qquad C_2 = \frac{\Gamma - \nu}{4} \tag{97}$$

The variance $\Sigma_t$ for the perturbation kernel $p(\mathbf{z}_t)\mathbf{z}_0$ is given by:

$$\boldsymbol{\Sigma}_t = \left( \begin{pmatrix} \Sigma_t^{xx} & \Sigma_t^{xm} \\ \Sigma_t^{xm} & \Sigma_t^{mm} \end{pmatrix} e^{-\left(\frac{\Gamma+\nu}{2}\right)\mathcal{B}(t)} \right) \otimes \boldsymbol{I}_d \tag{98}$$

where,

$$\Sigma_t^{xx} = A_1 \mathcal{B}^2(t)\Sigma_0^{xx} + A_2 \mathcal{B}^2(t)\Sigma_0^{mm} + A_3 \mathcal{B}(t)\Sigma_0^{xx} + A_4 \mathcal{B}^2(t) + A_5 \mathcal{B}(t) + (e^{2\lambda\mathcal{B}(t)} - 1) + \Sigma_0^{xx} \tag{99}$$

$$\Sigma_t^{xm} = C_1 \mathcal{B}^2(t)\Sigma_0^{xx} + C_2 \mathcal{B}^2(t)\Sigma_0^{mm} + C_3 \mathcal{B}(t)\Sigma_0^{xx} + C_4 \mathcal{B}(t)\Sigma_0^{mm} + C_5 \mathcal{B}^2(t) \tag{100}$$

$$\Sigma_t^{mm} = D_1 \mathcal{B}^2(t)\Sigma_0^{xx} + D_2 \mathcal{B}^2(t)\Sigma_0^{mm} + D_3 \mathcal{B}(t)\Sigma_0^{mm} + D_4 \mathcal{B}^2(t) + D_5 \mathcal{B}(t) + M(e^{2\lambda B(t)} - 1) + \Sigma_0^{mm} \tag{101}$$

where $\boldsymbol{\Sigma}_0 = \begin{pmatrix} \Sigma_0^{xx} & 0 \\ 0 & \Sigma_0^{mm} \end{pmatrix}$, $\mathcal{B}(t) = \int_0^t \beta(s)ds$ and coefficients:

$$A_1 = \frac{M^{-1}}{4} \qquad A_2 = \frac{M^{-2}}{4} \qquad A_3 = \frac{\nu - \Gamma}{2} \qquad A_4 = \frac{-M^{-1}}{2} \qquad A_5 = \frac{\Gamma - \nu}{2} \tag{102}$$

$$C_1 = \frac{\Gamma - \nu}{8} \qquad C_2 = \frac{(\Gamma - \nu)^3}{32} \qquad C_3 = \frac{-1}{2} \qquad C_4 = \frac{M^{-1}}{2} \qquad C_5 = \frac{\nu - \Gamma}{4} \tag{103}$$

$$D_1 = \frac{1}{4} \qquad D_2 = \frac{M^{-1}}{4} \qquad D_3 = \frac{\Gamma - \nu}{2} \qquad D_4 = \frac{-1}{2} \qquad D_5 = \frac{M(\nu - \Gamma)}{2} \tag{104}$$

It is worth noting that, when $\Gamma = 0$, $\bar{\nu} = M\nu$ and $\bar{\beta}(t) = \frac{\beta(t)}{2}$ such that $\mathcal{B}(t) = 2\bar{\mathcal{B}}(t)$ where $\bar{\mathcal{B}}(t) = \int_0^t \bar{\beta}(s)ds$, we have the following form of the mean $\mu_t$:

$$\mu_t = \begin{pmatrix} 2\bar{\nu}^{-2}\bar{\mathcal{B}}(t)\boldsymbol{x}_0 + 4\bar{\nu}^{-2}\bar{\mathcal{B}}(t)\boldsymbol{m}_0 + \boldsymbol{x}_0 \\ -\bar{\mathcal{B}}(t)\boldsymbol{x}_0 - 2\bar{\nu}^{-1}\bar{\mathcal{B}}(t)\boldsymbol{m}_0 + \boldsymbol{m}_0 \end{pmatrix} e^{-2\bar{\nu}^{-1}\bar{\mathcal{B}}(t)} \tag{105}$$

The expression for $\mu_t$ in Eqn. 105 is exactly the same as the mean of the perturbation kernel for CLD (Refer to Appendix B.1 in Dockhorn et al. [14]). A similar analysis holds for the variance $\Sigma_t$ which provides more insight into CLD being a special case of PSLD. Similar to CLD, at $t = 0$, we have $\boldsymbol{\mu}_0 = \begin{pmatrix} \mathbf{x}_0 \\ \mathbf{m}_0 \end{pmatrix}$, where $\mathbf{x}_0 \sim p(\mathbf{x}_0)$ (a.k.a the data generating distribution) and $\mathbf{m}_0 \sim \mathcal{N}(\mathbf{0}_d, M\gamma \boldsymbol{I}_d)$, where $\gamma$ is a scalar hyperparameter. Similarly, for DSM training, both $\Sigma_0^{xx}$ and $\Sigma_0^{mm}$ can be set to 0 (since $\mathbf{z}_t = [\mathbf{x}_t, \mathbf{m}_t]^T$ is a sample based estimate).

### B.3.2 Mean and Variance of $p(\mathbf{z}_t|\mathbf{x}_0)$

Since the data generating distribution for $m_0$ and the DSM perturbation kernel $p(\mathbf{z}_t|\mathbf{z}_0)$ are multivariate Gaussians, we can marginalize out the initial momentum variables $\mathbf{m}_0$ from $p(\mathbf{z}_t|\mathbf{z}_0)$ to obtain the perturbation kernel for HSM as $p(\mathbf{z}_t|\mathbf{x}_0) = \int p(\mathbf{z}_t|\mathbf{x}_0, \mathbf{m}_0)p(\mathbf{m}_0)d\mathbf{m}_0$. Consequently, the perturbation kernel $p(\mathbf{z}_t|\mathbf{x}_0)$ can be obtained by setting $m_0 = \mathbf{0}_d$ and $\Sigma_0^{mm} = M\gamma$ in the expressions of $\boldsymbol{\mu}_t$ and $\boldsymbol{\Sigma}_t$ for $p(\mathbf{z}_t|\mathbf{z}_0)$.

### B.3.3 Convergence

As $t \to \infty$, the mean $\mu_t$ converges to $\mathbf{0}_{2d}$ since the multiplicative term $e^{-(\frac{\nu+\Gamma}{4})\mathcal{B}(t)}$ goes to 0. Similarly, the covariance for the perturbation kernel converges to the following case:

$$\Sigma_e^{xx} = \lim_{t \to \infty} \Sigma_t^{xx} e^{-(\frac{\Gamma+\nu}{2})\mathcal{B}(t)} = 1 \tag{106}$$

$$\Sigma_e^{xm} = \lim_{t \to \infty} \Sigma_t^{xm} e^{-(\frac{\Gamma+\nu}{2})\mathcal{B}(t)} = 0 \tag{107}$$

$$\Sigma_e^{mm} = \lim_{t \to \infty} \Sigma_t^{mm} e^{-(\frac{\Gamma+\nu}{2})\mathcal{B}(t)} = M \tag{108}$$

Therefore, the perturbation kernel converges to the following steady-state distribution $p_{\text{EQ}}(\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{0}_d, \boldsymbol{I}_d)\mathcal{N}(\mathbf{m}; \mathbf{0}_d, M\boldsymbol{I}_d)$. It is worth noting that this is the exact equilibrium distribution that we specified in our SGM recipe to construct the forward process for PSLD.

### B.4 PSLD Sampling

The reverse SDE analogous to the forward SDE defined in Eqn. 54 can be formulated as follows [4]:

$$d\bar{\mathbf{z}}_t = \bar{\boldsymbol{f}}(\bar{\mathbf{z}}_t)dt + \boldsymbol{G}(T-t)d\bar{\mathbf{w}}_t \tag{109}$$

$$\bar{\boldsymbol{f}}(\bar{\mathbf{z}}_t) = \frac{\beta_t}{2} \begin{pmatrix} \Gamma\bar{\mathbf{x}}_t - M^{-1}\bar{\mathbf{m}}_t + 2\Gamma\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{0:d} \\ \bar{\mathbf{x}}_t + \nu\bar{\mathbf{m}}_t + 2M\nu\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{d:2d} \end{pmatrix}, \qquad \boldsymbol{G}(T-t) = \begin{pmatrix} \sqrt{\Gamma\beta_t} & 0 \\ 0 & \sqrt{M\nu\beta_t} \end{pmatrix} \otimes \boldsymbol{I}_d \tag{110}$$

where $\bar{\mathbf{z}}_t = \mathbf{z}_{T-t}$, $\bar{\mathbf{x}}_t = \mathbf{x}_{T-t}$, $\bar{\mathbf{m}}_t = \mathbf{m}_{T-t}$. Given an estimate of the score $\boldsymbol{s}_\theta(\mathbf{z}_t, T-t)$, one can simulate the above SDE to generate data from noise. Given $\bar{\mathbf{z}}_0 = (\bar{\mathbf{x}}_0, \bar{\mathbf{m}}_0)^T \sim p_{\text{EQ}}(\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{0}_d, \boldsymbol{I}_d)\mathcal{N}(\mathbf{m}; \mathbf{0}_d, M\boldsymbol{I}_d)$, we now discuss update steps for different samplers in context of PSLD.

### B.4.1 Euler-Maruyama (EM) Sampler

The EM update step for the reverse SDE corresponding to PSLD are as follows:

$$\begin{pmatrix} \bar{\mathbf{x}}_{t'} \\ \bar{\mathbf{m}}_{t'} \end{pmatrix} = \begin{pmatrix} \bar{\mathbf{x}}_t \\ \bar{\mathbf{m}}_t \end{pmatrix} + \frac{\beta_t \delta t}{2} \begin{pmatrix} \Gamma \bar{\mathbf{x}}_t - M^{-1} \bar{\mathbf{m}}_t + 2\Gamma s_\theta(\bar{\mathbf{z}}_t, T-t)|_{0:d} \\ \bar{\mathbf{x}}_t + \nu \bar{\mathbf{m}}_t + 2M\nu s_\theta(\bar{\mathbf{z}}_t, T-t)|_{d:2d} \end{pmatrix} + \begin{pmatrix} \sqrt{\Gamma \beta_t \delta t} \epsilon_{t'}^x \\ \sqrt{M\nu \beta_t \delta t} \epsilon_{t'}^m \end{pmatrix} \quad (111)$$

where $\epsilon_t = [\epsilon_t^x, \epsilon_t^m]^T \sim \mathcal{N}(\mathbf{0}_{2d}, \mathbf{I}_{2d})$ and $t' = t + \delta t$ where $\delta t$ is the step size for a single update.

### B.4.2 Symmetric Splitting CLD Sampler (SSCS)

Inspired by the application of splitting-based integrators in molecular dynamics [55], Dockhorn et al. [14] proposed the SSCS sampler with the following symmetric splitting scheme:

$$\begin{pmatrix} d\bar{\mathbf{x}}_t \\ d\bar{\mathbf{m}}_t \end{pmatrix} = \underbrace{\frac{\beta_t}{2} \begin{pmatrix} -\Gamma \bar{\mathbf{x}}_t - M^{-1} \bar{\mathbf{m}}_t \\ \bar{\mathbf{x}}_t - \nu \bar{\mathbf{m}}_t \end{pmatrix} dt + \mathbf{G}(T-t) d\bar{\mathbf{w}}_t}_{A} + \underbrace{\beta_t \begin{pmatrix} \Gamma \bar{\mathbf{x}}_t + \Gamma s_\theta(\bar{\mathbf{z}}_t, T-t)|_{0:d} \\ \nu \bar{\mathbf{m}}_t + M\nu s_\theta(\bar{\mathbf{z}}_t, T-t)|_{d:2d} \end{pmatrix} dt}_{S}$$

$$(112)$$

Dockhorn et al. [14] then approximate the flow map for the original SDE by the application of the following symmetric splitting schedule [69, 70]:

$$e^{t(\mathcal{L}_A + \mathcal{L}_S)} \approx \left[ e^{\frac{\delta t}{2} \mathcal{L}_A^*} e^{\delta t \mathcal{L}_S^*} e^{\frac{\delta t}{2} \mathcal{L}_A^*} \right]^N + \mathcal{O}(N\delta t^3) \quad (113)$$

where $N = \frac{t}{\delta t}$. The solution $\bar{\mathbf{z}}_t$ for the reverse SDE at any time t can then be obtained by the application of the flow map approximation of $e^{t(\mathcal{L}_A + \mathcal{L}_S)}$ to $\bar{\mathbf{z}}_0$. Since we use the same splitting formulation as Dockhorn et al. [14], the modified SSCS sampler for PSLD is still a first-order integrator sampler (Also see Appendix D in Dockhorn et al. [14] for more analysis of the SSCS sampler as proposed for CLD). However, since the form of the analytical splitting component in Eqn. 112 is different from CLD (due to a non-zero $\Gamma$), we next discuss the solution for this analytical form.

**Analytical splitting-term:** We have the following analytical splitting term:

$$\begin{pmatrix} d\bar{\mathbf{x}}_t \\ d\bar{\mathbf{m}}_t \end{pmatrix} = \frac{\beta_t}{2} \begin{pmatrix} -\Gamma \bar{\mathbf{x}}_t - M^{-1} \bar{\mathbf{m}}_t \\ \bar{\mathbf{x}}_t - \nu \bar{\mathbf{m}}_t \end{pmatrix} dt + \left( \begin{pmatrix} \sqrt{\Gamma \beta_t} & 0 \\ 0 & \sqrt{M\nu \beta_t} \end{pmatrix} \otimes \mathbf{I}_d \right) d\bar{\mathbf{w}}_t \quad (114)$$

The solution for this analytical SDE is similar to the derivation of the perturbation kernel in Appendix B.3. However, there are two key differences. Firstly, we need to integrate between time-intervals $(t, t + \delta t)$ as opposed to from $(0, t)$ for the perturbation kernel. Secondly, since we are sampling, we set the initial covariances $\Sigma_{xx}^t$ and $\Sigma_{mm}^t$ to zero. The analytical solution for the SDE in Eqn. 114 can then be specified as follows:

$$\bar{\mathbf{z}}_t \sim \mathcal{N}(\bar{\boldsymbol{\mu}}_t, \bar{\boldsymbol{\Sigma}}_t) \quad (115)$$

where

$$\boldsymbol{\mu}_{(\bar{\mathbf{x}}_t, \bar{\mathbf{m}}_t, t, t')} = \begin{pmatrix} A_1 \mathcal{B}(t, t') \bar{\boldsymbol{x}}_t + A_2 \mathcal{B}(t, t') \bar{\boldsymbol{m}}_t + \bar{\boldsymbol{x}}_t \\ C_1 \mathcal{B}(t, t') \bar{\boldsymbol{x}}_t + C_2 \mathcal{B}(t, t') \bar{\boldsymbol{m}}_t + \bar{\boldsymbol{m}}_t \end{pmatrix} e^{-(\frac{\nu+\Gamma}{4})\mathcal{B}(t,t')} \quad (116)$$

$$A_1 = \frac{\nu - \Gamma}{4} \qquad A_2 = \frac{-(\Gamma - \nu)^2}{8} \quad (117)$$

$$C_1 = \frac{1}{2} \qquad C_2 = \frac{\Gamma - \nu}{4} \quad (118)$$

The solution for the covariance is given by the following expression:

$$\boldsymbol{\Sigma}(t, t') = \left( \begin{pmatrix} \Sigma_{t,t'}^{xx} & \Sigma_{t,t'}^{xm} \\ \Sigma_{t,t'}^{xm} & \Sigma_{t,t'}^{mm} \end{pmatrix} e^{-(\frac{\Gamma+\nu}{2})\mathcal{B}(t,t')} \right) \otimes \mathbf{I}_d \quad (119)$$

where,

$$\Sigma_t^{xx} = -\frac{(\Gamma - \nu)^2}{8} \mathcal{B}^2(t, t') + \frac{(\Gamma - \nu)}{2} \mathcal{B}(t, t') + (e^{-(\frac{\Gamma+\nu}{2})\mathcal{B}(t,t')} - 1) \quad (120)$$

$$\Sigma_t^{xm} = \frac{(\Gamma - \nu)}{4} \mathcal{B}^2(t, t') \quad (121)$$

$$\Sigma_t^{mm} = -\frac{1}{2} \mathcal{B}^2(t, t') + \frac{M(\Gamma - \nu)}{2} \mathcal{B}(t, t') + M(e^{-(\frac{\Gamma+\nu}{2})\mathcal{B}(t,t')} - 1) \quad (122)$$

28

where $\mathcal{B}(t,t') = -\int_t^{t'} \beta(s)ds$ and $t' = t + \delta t$. Indeed, setting $\Gamma = 0$ recovers the original SSCS algorithm proposed in Dockhorn et al. [14]. Therefore, given $\bar{\mathbf{z}}_t = \left(\bar{\mathbf{x}}_t, \bar{\mathbf{m}}_t\right)^T$, the flow map update for the analytical splitting term $e^{\frac{\delta t}{2}\mathcal{L}_A^*}$ is given by:

$$\bar{\mathbf{z}}_{t'} \sim \mathcal{N}(\boldsymbol{\mu}(\bar{\mathbf{x}}_t, \bar{\mathbf{m}}_t, t, t'), \boldsymbol{\Sigma}(t, t')) \tag{123}$$

where $t' = t + \frac{\delta t}{2}$.

**Score-based splitting term:** The flow map update for the score-based splitting term $e^{\delta t\mathcal{L}_S^*}$ is given by an Euler update as follows:

$$\begin{pmatrix} \bar{\mathbf{x}}_{t'} \\ \bar{\mathbf{m}}_{t'} \end{pmatrix} = \begin{pmatrix} \bar{\mathbf{x}}_t \\ \bar{\mathbf{m}}_t \end{pmatrix} + \delta t \beta_t \begin{pmatrix} \Gamma\bar{\mathbf{x}}_t + \Gamma\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{0:d} \\ \nu\bar{\mathbf{m}}_t + M\nu\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{d:2d} \end{pmatrix} \tag{124}$$

Combining the two splitting terms together, a more generic form of the SSCS algorithm can be specified as follows:

---

**Algorithm 1** *Modified SSCS* (Terms in blue indicate differences from the SSCS sampler proposed in Dockhorn et al. [14])

---

**Input:** Trajectory length T, Score function $\boldsymbol{s}_\theta(\mathbf{z}_t, T-t)$, PSLD parameters $\Gamma$, $\nu$, $\beta_t$, $M = \frac{(\Gamma-\nu)^2}{4}$, number of sampling steps $N$, step sizes $\{\delta t_n \geq 0\}_{n=0}^{N-1}$ spanning the interval $(0, T-\epsilon)$.
**Output:** $\bar{\mathbf{z}}_T = (\bar{\mathbf{x}}_T, \bar{\mathbf{m}}_T)$

$\bar{\mathbf{x}}_0 \sim \mathcal{N}(\mathbf{0}_d, \boldsymbol{I}_d), \bar{\mathbf{m}}_0 \sim \mathcal{N}(\mathbf{0}_d, M\boldsymbol{I}_d), \bar{\mathbf{z}}_0 = (\bar{\mathbf{x}}_0, \bar{\mathbf{m}}_0)$      ▷ Draw initial prior samples from $p_{\text{EQ}}(\mathbf{u})$
$t = 0$                         ▷ Initialize time
**for** $n = 0$ **to** $N-1$ **do**
 $\bar{\mathbf{z}}_{n+\frac{1}{2}} \sim \mathcal{N}(\boldsymbol{\mu}(\bar{\mathbf{x}}_n, \bar{\mathbf{m}}_n, t, t+\frac{\delta t_n}{2}), \boldsymbol{\Sigma}(t, t+\frac{\delta t_n}{2}))$      ▷ First half-step: Apply $\exp\{\frac{\delta t_n}{2}\hat{\mathcal{L}}_A^*\}$
 $\bar{\mathbf{z}}_{n+\frac{1}{2}} \leftarrow \bar{\mathbf{z}}_{n+\frac{1}{2}} + \delta t_n\beta_t \begin{pmatrix} \Gamma\bar{\mathbf{x}}_{n+\frac{1}{2}} + \Gamma\boldsymbol{s}_\theta(\bar{\mathbf{z}}_{n+\frac{1}{2}}, T-t)|_{0:d} \\ \nu\bar{\mathbf{m}}_{n+\frac{1}{2}} + M\nu\boldsymbol{s}_\theta(\bar{\mathbf{z}}_{n+\frac{1}{2}}, T-t)|_{d:2d} \end{pmatrix}$  ▷ Full step: Apply $\exp\{\delta t_n\hat{\mathcal{L}}_S^*\}$
 $\bar{\mathbf{z}}_{n+1} \sim \mathcal{N}(\boldsymbol{\mu}(\bar{\mathbf{x}}_{n+\frac{1}{2}}, \bar{\mathbf{m}}_{n+\frac{1}{2}}, t, t+\frac{\delta t_n}{2}), \boldsymbol{\Sigma}(t, t+\frac{\delta t_n}{2}))$  ▷ Second half-step: Apply $\exp\{\frac{\delta t_n}{2}\hat{\mathcal{L}}_A^*\}$
 $t \leftarrow t + \delta t_n$                    ▷ Update time
**end for**
$\bar{\mathbf{z}}_N \leftarrow \bar{\mathbf{z}}_N + \epsilon\frac{\beta_t}{2} \begin{pmatrix} \Gamma\bar{\mathbf{x}}_{n+1} - M^{-1}\bar{\mathbf{m}}_{n+1} + 2\Gamma\boldsymbol{s}_\theta(\bar{\mathbf{z}}_{n+1}, \epsilon)|_{0:d} \\ \bar{\mathbf{x}}_{n+1} + \nu\bar{\mathbf{m}}_{n+1} + 2M\nu\boldsymbol{s}_\theta(\bar{\mathbf{z}}_{n+1}, \epsilon)|_{d:2d} \end{pmatrix}$     ▷ Denoising
$(\bar{\mathbf{x}}_N, \bar{\mathbf{m}}_N) = \bar{\mathbf{z}}_N$         ▷ Extract output data and momentum samples

---

### B.4.3 Probability Flow ODE

Following Song et al. [4], the probability flow ODE for PSLD can be specified as follows:

$$d\bar{\mathbf{z}}_t = \bar{\boldsymbol{f}}(\bar{\mathbf{z}}_t)dt \tag{125}$$

$$\bar{\boldsymbol{f}}(\bar{\mathbf{z}}_t) = \frac{\beta_t}{2} \begin{pmatrix} \Gamma\bar{\mathbf{x}}_t - M^{-1}\bar{\mathbf{m}}_t + \Gamma\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{0:d} \\ \bar{\mathbf{x}}_t + \nu\bar{\mathbf{m}}_t + M\nu\boldsymbol{s}_\theta(\bar{\mathbf{z}}_t, T-t)|_{d:2d} \end{pmatrix} \tag{126}$$

The Probability-Flow ODE can be solved using any fixed/adaptive step-size black-box ODE solvers like RK45 [43]

## C Implementation Details

### C.1 Datasets and Preprocessing

We use CIFAR-10 [18] (50k images) and CelebA-64 ($\approx$ 200k images) [19] datasets for both quantitative and qualitative analysis. We use the AFHQv2 [33] dataset ($\approx$ 14k images) only for qualitative analysis. Unless specified otherwise, we always use the CelebA dataset at 64x64 resolution and the AFHQv2 dataset at 128 x 128 resolution. During training, all datasets are preprocessed to a numerical range of [-1, 1]. Following prior work, we use random horizontal flips to train all models (ablation and SOTA) across datasets as a data augmentation strategy.

| Hyperparameter | CIFAR-10 | | CelebA-64 | | AFHQv2 |
| --- | --- | --- | --- | --- | --- |
| | SOTA | Ablation | SOTA | Ablation | Qualitative |
| Base channels | 128 | 128 | 128 | 128 | 128 |
| Channel multiplier | [2,2,2] | [1,2,2,2] | [1,2,2,2] | [1,1,2,2,2] | [1,2,2,2,3] |
| # Residual blocks | 4,8 | 2 | 4 | 4 | 2 |
| Non-Linearity | Swish | Swish | Swish | Swish | Swish |
| Attention resolution | [16] | [16] | [16] | [16] | [16] |
| # Attention heads | 1 | 1 | 1 | 1 | 1 |
| Dropout | 0.15 | 0.1 | 0.1 | 0.1 | 0.2 |
| Finite Impulse Response (FIR) [72] | True | False | True | False | False |
| FIR kernel | [1,3,3,1] | N/A | [1,3,31] | N/A | N/A |
| Progressive Input | Residual | None | Residual | None | None |
| Progressive Combine | Sum | Sum | Sum | Sum | Sum |
| Embedding type | Fourier | Positional | Fourier | Positional | Positional |
| Sigma scaling | False | False | False | False | False |
| Model size | 55M/97M | 39M | 62M | 66M | 68M |

Table 8: Score Network hyperparameters for PSLD.

## C.2 Score Network Architecture

Table 8 illustrates our score model architectures for different datasets. Our network architectures are largely based on the design of the DDPM++/NCSN++ score networks introduced in Song et al. [4]. Apart from minor design choices, the DDPM++/NCSN++ score-network architectures are primarily based on the U-Net [71] model. We further highlight several key aspects of our score network architectures across different datasets as follows:

**CIFAR-10**: We use a smaller version (39M) of the DDPM++ architecture (with the number of residual blocks per resolution set to two) for ablation studies (for both VP-SDE and PSLD) while we use the NCSN++ architecture [4] for training larger models used for SOTA comparisons. Moreover, when training larger models, like Karras et al. [31] we remove the layers at 4x4 resolution and re-distribute capacity to the layers at the 16x16 resolution. This results in model sizes of 55M/97M parameters corresponding to four and eight residual blocks per resolution, respectively, with channel multipliers [2,2,2]. Moreover, when training larger models (55M/97M), we slightly increase the dropout rate from 0.1 to 0.15. We observed that these changes improved performance slightly while reducing model sizes and enabling faster training.

**CelebA-64**: Similar to CIFAR-10, we use a DDPM++ score model architecture for ablation experiments. while we use a NCSN++ architecture for SOTA comparisons. Moreover, we remove the 4x4 layers from our ablation model for SOTA analysis and increase the channel multiplier for the 32x32 layers from 1 to 2. The dropout rate is set to 0.1 due to a larger dataset size for CelebA-64. This setting results in a model size of approximately 66M for the ablation experiments and 62M for SOTA comparisons.

**AFHQv2**: We use the original DDPM++ architecture for training our AFHQv2 model. Additionally, we increase the dropout rate to 0.2, given a relatively smaller dataset size. This setting results in a model size of approximately 68M parameters for qualitative analysis.

## C.3 SDE Parameters

**PSLD:** For PSLD (including the CLD baseline), unless specified otherwise, we set the mass parameter $M^{-1} = 4$ and $\beta = 8$. The choice of these parameters is motivated by empirical results presented in Dockhorn et al. [14]. We add a stabilizing numerical epsilon value of $1e^{-9}$ in the diagonal entries of the Cholesky decomposition of $\Sigma_t$ when sampling the perturbed data-point $\mathbf{z}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)$ during training. The data-generating distribution is set to $p_0(\mathbf{z}) = \mathcal{N}(\mathbf{0}_d, \boldsymbol{I}_d)\mathcal{N}(\mathbf{0}_d, M\gamma\boldsymbol{I}_d)$ where $\gamma = 0.04$. For SOTA analysis, we experiment with $\Gamma \in \{0.01, 0.02\}$ for CIFAR-10 and $\Gamma = 0.005$ for the CelebA-64 datasets. We chose these values of $\Gamma$ and $\nu$ based on the best-performing (in terms of FID)

|  | CIFAR-10 | | CelebA-64 | | AFHQv2 |
|---|---|---|---|---|---|
|  | SOTA | Ablation | SOTA | Ablation | Qualitative |
| Random Seed | 0 | 0 | 0 | 0 | 0 |
| # iterations | 800k | 800k | 800k | 320k | 400k |
| Optimizer | Adam | Adam | Adam | Adam | Adam |
| Grad Clip. cutoff | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| Learning rate (LR) | 2e-4 | 2e-4 | 2e-4 | 2e-4 | 1e-4 |
| LR Warmup steps | 5000 | 5000 | 5000 | 5000 | 5000 |
| FP16 | False | False | False | False | False |
| EMA Rate | 0.9999 | 0.9999 | 0.9999 | 0.9999 | 0.9999 |
| Effective Batch size | 128 | 128 | 128 | 128 | 64 |
| # GPUs | 8 | 4 | 8 | 4 | 8 |
| Train eps cutoff | 1e-5 | 1e-5 | 1e-5 | 1e-5 | 1e-5 |

Table 9: Training hyperparameters for PSLD

|  |  | CIFAR-10 | AFHQ-v2 |
|---|---|---|---|
|  | Random Seed | 0 | 0 |
|  | # iterations | 200k | 70k |
|  | Optimizer | Adam | Adam |
|  | LR | 2e-4 | 2e-4 |
| Training | Warmup steps | 5000 | 5000 |
|  | FP16 | False | False |
|  | Batch size | 256 | 64 |
|  | # GPUs | 4 | 4 |
|  | Train eps cutoff | 1e-5 | 1e-5 |
| SDE | $M^{-1}$ | 4.0 | 4.0 |
|  | $\Gamma$ | 0.01 | 0.01 |
|  | $\nu$ | 4.01 | 4.01 |
|  | $\beta$ | 8.0 | 8.0 |

Table 10: Classifier Training Hyperparameters

|  | CIFAR-10 | AFHQ-v2 |
|---|---|---|
| Base channels | 128 | 128 |
| Num. classes | 10 | 3 |
| Channel multiplier | [1,2,3,4] | [1,2,3,4] |
| # Residual blocks | 4 | 4 |
| Non-Linearity | Swish | Swish |
| Attention resolution | [16, 8] | [16, 8] |
| # Attention heads | 1 | 1 |
| Dropout | 0.1 | 0.1 |
| FIR [72] | False | False |
| Progressive Input | None | None |
| Progressive Combine | Sum | Sum |
| Embedding type | Positional | Positional |
| Sigma scaling | False | False |
| Model size | 56.7M | 57.8M |

Table 11: Classifier Network Hyperparameters

ablation models for these datasets (See Table 4 in the main text). Lastly, for training our AFHQv2 model for qualitative analysis, we set $\Gamma = 0.01$. All the other SDE parameters remain the same.

**VP-SDE:** For our VP-SDE baseline, following Song et al. [4], we set $\beta_{\min} = 0.1$ and $\beta_{\max} = 20.0$

## C.4 Training

Table 9 summarizes the different training hyperparameters across datasets and evaluation settings (ablation and SOTA). Additionally, we use the Hybrid Score Matching (HSM) objective (See Appendix B.2.1) for all augmented state-space models (PSLD and CLD); for the VP-SDE baseline, we use the Denoising Score Matching (DSM) objective. Throughout this work, we optimize for sample quality and thus use the *epsilon-prediction* loss during training (See Appendix B.2.1).

## C.5 Evaluation

**SDE Sampling:** As is common in prior works [4, 14], we use the integration interval $(1e^{-3}, 1.0)$ for solving the reverse SDE/ODE for sample generation. Unless specified otherwise, we use 1000 sampling steps when using numerical SDE solvers. When using numerical black-box ODE solvers, we use the RK-45 [43] solver with the same absolute and relative tolerance levels. We use the ODE solver at different tolerance levels for ablations studies (Table 6 in the main text) and tolerance levels of $1e^{-5}$ and $1e^{-4}$ for reporting SOTA results on CIFAR-10. Similarly, we use a tolerance level of $1e^{-5}$ for reporting ODE solver performance on CelebA-64. We use the odeint function from the torchdiffeq[73] package when using the black-box ODE solver for sampling.

**Timestep Selection during Sampling**: We use *Uniform* (US) and *Quadratic* (QS) striding for timestep discretization in this work. In uniform striding, given an NFE budget N, we discretize the integration interval ($\epsilon$, T) into N equidistant parts, which are then used for score function evaluations. In quadratic striding [14, 39], the evaluation timepoints are given by:

$$\tau_i = \left(\frac{i}{N}\right)^2 \quad \forall i \in [0, N) \tag{127}$$

This ensures more number of score function evaluations in the lower timestep regime (i.e. $t$, which is close to the data). This kind of timestep selection is particularly useful when the NFE budget is limited (See Table 5).

**Last-Step Denoising**: Similar to prior works [4, 14, 45], we perform a single denoising EM step (without noise injection) at the very last step of our sampling routine for both SDE and ODE solvers. Formally, we perform the following update:

$$\begin{pmatrix} \mathbf{x}_0 \\ \mathbf{m}_0 \end{pmatrix} = \begin{pmatrix} \mathbf{x}_\epsilon \\ \mathbf{m}_\epsilon \end{pmatrix} + \frac{\beta_t \epsilon}{2} \begin{pmatrix} \Gamma \mathbf{x}_\epsilon - M^{-1}\mathbf{m}_\epsilon + 2\Gamma s_\theta(\mathbf{z}_\epsilon, \epsilon)|_{0:d} \\ \mathbf{x}_\epsilon + \nu\mathbf{m}_\epsilon + 2M\nu s_\theta(\mathbf{z}_\epsilon, \epsilon)|_{d:2d} \end{pmatrix} \tag{128}$$

where $\epsilon = 1e^{-3}$. Such a denoising step has been found useful in removing additional noise, thereby improving FID scores [45].

**Evaluation Metrics:** For most ablation experiments involving the analysis of speed vs sample quality trade-offs between different models, we report the FID [20] score on 10k samples for computational convenience. For SOTA comparisons, we report FID for 50k samples for both CIFAR-10 and CelebA-64 datasets. When reporting extended SOTA results for CIFAR-10, we also report the Inception Score (IS) [21] metric. We use the `torch-fidelity`[74] package for computing all FID and IS scores reported in this work. When reporting average NFE (number of function evaluations) in Table 6, we average the NFE values over a batch size of 16 samples for 10k samples in aggregate and take a ceiling of the resulting value.

## C.6 Classifier Architecture and Training

For class conditional synthesis (Appendix D.4), we append the downsampling part of the UNet architecture with a classification head and use the resulting model as our classifier architecture. Table 11 shows different hyperparameters of our classifier model architecture. For classifier training, we set $\Gamma = 0.01$ for the AFHQv2 and the CIFAR-10 datasets. The remaining SDE parameters remain unchanged from our previous setting. Table 10 lists different hyperparameters for classifier training.

# D Additional Results

## D.1 Impact of $\Gamma$ and $\nu$ on PSLD Sample Quality

Table 12 shows the impact of varying $\Gamma$ and $\nu$ (with a fixed $M^{-1}$) on the PSLD sample quality using the EM-sampler with quadratic striding (EM-QS) for the CIFAR-10 dataset. Extending Table 4, we additionally present FID scores for $\Gamma \in \{4.0, 8.0\}$ in Table 12. As we increase the value of $\Gamma$ to 4.0 or 8.0, the FID scores further increase to 11.43 and 14.15, respectively, confirming our observations in Section 4.2. Figure 6 further illustrates the qualitative impact of increasing $\Gamma$ on CIFAR-10 sample quality. For the setting with $\Gamma = 8.0$, using EM with uniform striding (EM-US) introduces evident noise artifacts in the generated samples. However, such artifacts are less pronounced when using EM with quadratic striding ((EM-QS)) instead. This suggests potential denoising problems for low timestep indices during sampling as quadratic striding focuses more score network evaluations in the low timestep regime, which might lead to lesser artifacts. This is similar to our observations in Section 4.2 (See Figure 7 for more qualitative results on CelebA-64). We now provide a formal justification for these observations.

Given an input sample $\bar{\mathbf{z}}_t = (\bar{\mathbf{x}}_t, \bar{\mathbf{m}}_t)$ at time t, consider the following update rule for the EM-sampler for PSLD with a uniform spacing interval of $\delta t$ between successive steps:

$$\begin{pmatrix} \bar{\mathbf{x}}_{t'} \\ \bar{\mathbf{m}}_{t'} \end{pmatrix} = \begin{pmatrix} \bar{\mathbf{x}}_t \\ \bar{\mathbf{m}}_t \end{pmatrix} + \frac{\beta\delta t}{2} \begin{pmatrix} \Gamma\bar{\mathbf{x}}_t - M^{-1}\bar{\mathbf{m}}_t + 2\Gamma s_\theta(\bar{\mathbf{z}}_t, T-t)|_{0:d} \\ \bar{\mathbf{x}}_t + \nu\bar{\mathbf{m}}_t + 2M\nu s_\theta(\bar{\mathbf{z}}_t, T-t)|_{d:2d} \end{pmatrix} + \begin{pmatrix} \sqrt{\Gamma\beta\delta t}\boldsymbol{\epsilon}_{t'}^x \\ \sqrt{M\nu\beta\delta t}\boldsymbol{\epsilon}_{t'}^m \end{pmatrix} \tag{129}$$

| $\Gamma$ | $\nu$ | $M^{-1} = \frac{(\Gamma - \nu)^2}{4}$ | FID@50k ↓ (EM-QS) |
|---|---|---|---|
| 0 | 4 | 4 | 3.64 |
| 0.005 | 4.005 | 4 | 3.42 |
| 0.01 | 4.01 | 4 | 3.15 |
| 0.02 | 4.02 | 4 | 3.26 |
| 0.25 | 4.25 | 4 | 4.99 |
| 4 | 0 | 4 | 11.43 |
| 8 | 4 | 4 | 14.15 |

Table 12: Extended results for impact of the choice of $\Gamma$ on sample quality for CIFAR-10. FID (lower is better) reported on 50k samples.

To simplify notation, let us denote $\bar{\beta} = \frac{\beta \delta t}{2}$, $s_\theta^x(\bar{\mathbf{z}}_t) = s_\theta(\bar{\mathbf{z}}_t, T - t)|_{0:d}$ and $s_\theta^m(\bar{\mathbf{z}}_t) = s_\theta(\bar{\mathbf{z}}_t, T - t)|_{d:2d}$. Therefore, for the next timestep $t'$, we have:

$$\bar{\mathbf{x}}_{t'} = \bar{\mathbf{x}}_t + \bar{\beta}\Big(\Gamma \bar{\mathbf{x}}_t - M^{-1}\bar{\mathbf{m}}_t + 2\Gamma s_\theta^x(\bar{\mathbf{z}}_t)\Big) + \sqrt{\Gamma \beta \delta t}\, \epsilon_{t'}^x \tag{130}$$

$$\bar{\mathbf{m}}_{t'} = \bar{\mathbf{m}}_t + \bar{\beta}\Big(\bar{\mathbf{x}}_t + \nu \bar{\mathbf{m}}_t + 2M\nu s_\theta^m(\bar{\mathbf{z}}_t)\Big) + \sqrt{M\nu \beta \delta t}\, \epsilon_{t'}^m \tag{131}$$

Similarly for the next consecutive time-step $t''$, we have the following EM-update rule:

$$\bar{\mathbf{x}}_{t''} = \bar{\mathbf{x}}_{t'} + \bar{\beta}\Big(\Gamma \bar{\mathbf{x}}_{t'} - M^{-1}\bar{\mathbf{m}}_{t'} + 2\Gamma s_\theta^x(\bar{\mathbf{z}}_{t'})\Big) + \sqrt{\Gamma \beta \delta t}\, \epsilon_{t''}^x \tag{132}$$

Substituting the update expressions for $\bar{\mathbf{x}}_{t'}$ and $\bar{\mathbf{m}}_{t'}$ from Eqns. 130-131 in the update rule for $\bar{\mathbf{x}}_{t''}$, we have the following modified update rule for $\bar{\mathbf{x}}_{t''}$:

$$\bar{\mathbf{x}}_{t''} = \boldsymbol{f}(\bar{\mathbf{x}}_t, \bar{\mathbf{m}}_t) + \hat{s}_\theta + \boldsymbol{\eta} \tag{133}$$

where $\boldsymbol{f}$ is a function of $(\bar{\mathbf{x}}_t, \bar{\mathbf{m}}_t)$, $\boldsymbol{\eta}$ is the aggregate stochastic noise. More importantly, the score term $\hat{s}_\theta$ is given as follows:

$$\hat{s}_\theta = 2\bar{\beta}\Gamma s_\theta^x(\bar{\mathbf{z}}_t) + 2\bar{\beta}^2\Big[\Gamma^2 s_\theta^x(\bar{\mathbf{z}}_t) - \nu s_\theta^m(\bar{\mathbf{z}}_t)\Big] + 2\Gamma\bar{\beta}s_\theta^x(\bar{\mathbf{z}}_{t'}) \tag{134}$$

$$= 2\bar{\beta}\Gamma s_\theta^x(\bar{\mathbf{z}}_t) + 2\bar{\beta}^2\Bigg[\begin{pmatrix} s_\theta^x(\bar{\mathbf{z}}_t) \\ s_\theta^m(\bar{\mathbf{z}}_t) \end{pmatrix}^T \begin{pmatrix} \Gamma^2 \\ -\nu \end{pmatrix}\Bigg] + 2\Gamma\bar{\beta}s_\theta^x(\bar{\mathbf{z}}_{t'}) \tag{135}$$

$$= 2\bar{\beta}\Gamma s_\theta^x(\bar{\mathbf{z}}_t) + 2\bar{\beta}^2\Big[s_\theta(\bar{\mathbf{z}}_t)^T \begin{pmatrix} \Gamma^2 \\ -\nu \end{pmatrix}\Big] + 2\Gamma\bar{\beta}s_\theta^x(\bar{\mathbf{z}}_{t'}) \tag{136}$$

In this work, we parameterize the score $s_\theta(\bar{\mathbf{z}}_t) = -\boldsymbol{L}_t^{-T}\epsilon_\theta(\bar{\mathbf{z}}_t)$ where $\boldsymbol{L}_t$ is Cholesky factorization matrix of the covariance matrix $\boldsymbol{\Sigma}_t$ of the perturbation kernel at time t. Substituting this parameterization in Eqn. 136, we get,

$$\hat{s}_\theta = 2\bar{\beta}\Gamma s_\theta^x(\bar{\mathbf{z}}_t) + 2\bar{\beta}^2\Big[-\epsilon_\theta^T(\bar{\mathbf{z}}_t)\boldsymbol{L}_t^{-1}\begin{pmatrix} \Gamma^2 \\ -\nu \end{pmatrix}\Big] + 2\Gamma\bar{\beta}s_\theta^x(\bar{\mathbf{z}}_{t'}) \tag{137}$$

$$= 2\bar{\beta}\Gamma s_\theta^x(\bar{\mathbf{z}}_t) - 2\bar{\beta}^2\Big[\Gamma^2(l_t^{xx}\epsilon_\theta^x(\bar{\mathbf{z}}_t) + l_t^{xm}\epsilon_\theta^m(\bar{\mathbf{z}}_t)) - \nu l_t^{mm}\epsilon_\theta^m(\bar{\mathbf{z}}_t)\Big] + 2\Gamma\bar{\beta}s_\theta^x(\bar{\mathbf{z}}_{t'}) \tag{138}$$

$$= 2\bar{\beta}\Gamma s_\theta^x(\bar{\mathbf{z}}_t) - 2\bar{\beta}^2\Big[\underbrace{\Gamma^2 l_t^{xx}}_{=\lambda_1}\epsilon_\theta^x(\bar{\mathbf{z}}_t) + \underbrace{(\Gamma^2 l_t^{xm} - \nu l_t^{mm})}_{=\lambda_2}\epsilon_\theta^m(\bar{\mathbf{z}}_t)\Big] + 2\Gamma\bar{\beta}s_\theta^x(\bar{\mathbf{z}}_{t'}) \tag{139}$$

where,

$$l_t^{xx} = \frac{1}{\sqrt{\Sigma_t^{xx}}} \tag{140}$$

$$l_t^{xm} = \frac{-\Sigma_t^{xm}}{\sqrt{\Sigma_t^{xx}}\sqrt{\Sigma_t^{xx}\Sigma_t^{mm} - (\Sigma_t^{xm})^2}} \tag{141}$$
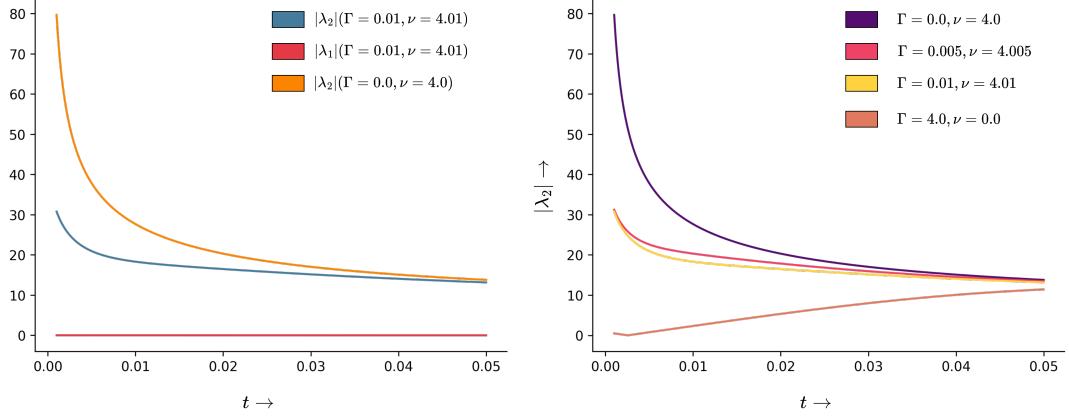
Figure 4: (a) Comparison between $|\lambda_1|$ and $|\lambda_2|$ corresponding to $\Gamma = 0.0$ and $\Gamma = 0.01$ in the low-timestep regime for a fixed $M^{-1} = 4$. (b) Variation of $|\lambda_2|$ for different values of $(\Gamma, \nu)$

$$l_t^{mm} = \sqrt{\frac{\Sigma_t^{xx}}{\Sigma_t^{xx}\Sigma_t^{mm} - \left(\Sigma_t^{xm}\right)^2}} \tag{142}$$

Assuming the input $\bar{\mathbf{z}}_t$ is a sample from the underlying flow map of the reverse SDE (a very strong assumption), the score term $\hat{s}_\theta$ in Eqn. 132 is the primary source of introducing errors (since the neural network-based score prediction will be offset by some error from the true underlying score). Without loss of generality, we further assume that the update timepoints $t, t'$ and $t''$ lie in the low timestep regime. Furthermore, for notational convenience, we denote $\lambda_1 = \Gamma^2 l_t^{xx}$ and $\lambda_2 = (\Gamma^2 l_t^{xm} - \nu l_t^{mm})$ as the scaling factors for the second term in Eqn. 139. Thus, the error introduced due to the neural network predictors $\epsilon_\theta^x(\bar{\mathbf{z}}_t)$ and $\epsilon_\theta^m(\bar{\mathbf{z}}_t)$ will be scaled by $\lambda_1$ and $\lambda_2$ respectively. Therefore, for achieving lower sampler discretization errors, it might be desirable to have low magnitudes of $\lambda_1$ and $\lambda_2$. We now qualitatively analyze the magnitude of these coefficients for different ranges of values of $\Gamma$ and $\nu$.

**Case-1: Effect of using a non-zero** $\Gamma$: We first analyze the impact of using a non-zero $\Gamma$ value on the magnitude of $\lambda_1$ and $\lambda_2$. Figure 4a illustrates the impact of the choice of $\Gamma$ and $\nu$ on the coefficients $\lambda_1$ and $\lambda_2$ in the low-timestep regime. When $\Gamma = 0$, the error in the score term $\hat{s}_\theta$ will be only due to the term $|\lambda_2|\epsilon_\theta^m(\bar{\mathbf{z}}_t)$. As illustrated in Figure 4a, the value of $|\lambda_2|$ (when $\Gamma = 0$) is very high in the low-timestep regime and, therefore might negatively impact the sample quality since any errors in the estimation of $\epsilon_\theta^m(\bar{\mathbf{z}}_t)$ would be scaled by a large factor.

Interestingly, for the setting $\Gamma = 0.01, \nu = 4.01$, the value of $|\lambda_2|$ reduces significantly, thus reducing the error scaling factor. It is worth noting that using a non-zero $\Gamma$ also simultaneously enables error contribution from other terms in $\hat{s}_\theta$ involving $\Gamma$ (especially $|\lambda_1|\epsilon_\theta^x(\bar{\mathbf{z}}_t)$). However, as illustrated in Figure 4a, the value of $|\lambda_1|$ is extremely small as compared to $|\lambda_2|$ making the additional error introduced insignificant. Due to this reason, the overall error introduced by the score term $\hat{s}_\theta$ is more when $\Gamma = 0$ as compared to the setting with a (small) non-zero $\Gamma$ value. This explains why using a small value of $\Gamma$ yields better sample quality than our CLD baseline (See Table 4)

**Case-2: Effect of using a large** $\Gamma$: Figure 4b illustrates the variation of $|\lambda_2|$ for some more values of $\Gamma$. Interestingly for $\Gamma = 4.0$, the value of $|\lambda_2|$ decreases to almost 0 in the low-timestep regime. However, for $\Gamma = 4.0$, the value of $|\lambda_1|$ increases significantly (Figure 5a), therefore, leading to large error scaling factors in the term $|\lambda_1|\epsilon_\theta^x(\bar{\mathbf{z}}_t)$. This finding justifies our observation in Figure 2, where a value of $\Gamma = 0.25$ makes sample quality significantly worse for the CelebA-64 dataset and is unable to recover high-frequency details. Figure 5b further illustrates the variation of $|\lambda_1|$ for different $(\Gamma, \nu)$ pairs in the low-timestep regime.

From the above analysis, it seems that *the choice of $\Gamma$ provides an important trade-off between balancing the errors produced due to the terms $\epsilon_\theta^x(\bar{\mathbf{z}}_t)$ and $\epsilon_\theta^m(\bar{\mathbf{z}}_t)$ in Eqn. 139.* Therefore, the choice of $\Gamma$ is crucial for sample quality in PSLD.
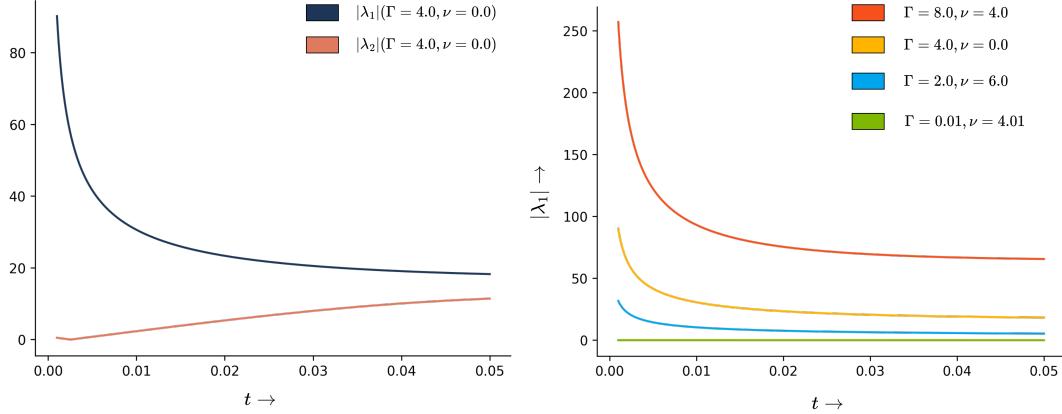
34

Figure 5: (a) Comparison between $|\lambda_1|$ and $|\lambda_2|$ corresponding to $\Gamma = 4.0$ and $\nu = 0.0$ in the low-timestep regime. (b) Variation of $|\lambda_1|$ for different values of $(\Gamma, \nu)$

| Sampler | Method | NFE (FID@10k ↓) | | | | |
|---|---|---|---|---|---|---|
| | | 50 | 100 | 250 | 500 | 1000 |
| EM-QS | CLD | 25.01 | 8.91 | 5.97 | 5.61 | 5.7 |
| | VP-SDE | **17.72** | 7.45 | 5.59 | 5.51 | 5.51 |
| | (Ours) PSLD ($\Gamma = 0.01$) | 23.96 | 8.12 | 5.41 | **5.13** | **5.24** |
| | (Ours) PSLD ($\Gamma = 0.02$) | 19.94 | **7.33** | **5.26** | 5.20 | 5.28 |
| EM-US | CLD | 119.68 | 45.60 | 9.08 | 5.71 | 5.65 |
| | VP-SDE | **84.54** | 41.93 | 12.61 | 5.92 | 5.19 |
| | (Ours) PSLD ($\Gamma = 0.01$) | 109.01 | 40.22 | **9.07** | 5.25 | 4.95 |
| | (Ours) PSLD ($\Gamma = 0.02$) | 100.62 | **39.96** | 11.26 | 5.45 | **4.82** |
| SSCS-QS | CLD | 21.31 | 8.37 | 5.82 | 5.75 | 5.69 |
| | (Ours) PSLD ($\Gamma = 0.01$) | 18.41 | 7.42 | 5.41 | **5.28** | 5.29 |
| | (Ours) PSLD ($\Gamma = 0.02$) | **16.12** | **7.16** | 5.36 | 5.35 | **5.27** |
| SSCS-US | CLD | 75.45 | 24.74 | 6.09 | 5.74 | 5.78 |
| | (Ours) PSLD ($\Gamma = 0.01$) | 76.6 | 21.25 | **5.18** | 5.10 | 5.33 |
| | (Ours) PSLD ($\Gamma = 0.02$) | **72.42** | **20.46** | 5.19 | **4.92** | **5.29** |

Table 13: Extended Speed vs. Sample quality comparisons using the SDE setup for CIFAR-10. FID computed for 10k samples. Values in **bold** indicate the best result for that column.

## D.2 Additional Speed vs. Sample Quality Comparisons

**CIFAR-10**: We extend the Speed vs. Sample Quality results in Table 5 to include results for PSLD with $\Gamma = 0.01$ in Table 13

**CelebA-64**: Similar to our setup for CIFAR-10 (Section 4.3), we benchmark the speed vs quality tradeoffs of PSLD ($\Gamma = 0.005$) against our CLD ablation baseline for the CelebA-64 dataset (See Table 14). Similar to CIFAR-10, PSLD outperforms our CLD baseline across all timesteps. The performance difference is most notable in the low-timestep regime (FID of 6.99 for PSLD with EM-QS vs 10.7 for CLD with SSCS-QS). However, there are two notable differences in our observations when compared to CIFAR-10:

**(i)** Firstly, quadratic striding works best for the CelebA-64 dataset. This contrasts with CIFAR-10, where a uniform striding schedule works better.

**(ii)** More interestingly, PSLD achieves the best performance of FID=3.77 at N=250 steps, and sample quality degrades on further increasing the number of steps. This contrasts our results for CIFAR-10, where PSLD achieves the best performance at T=1000.

| | | NFE (FID@10k $\downarrow$) | | | | |
|---|---|---|---|---|---|---|
| Sampler | Method | 50 | 100 | 250 | 500 | 1000 |
| EM-QS | CLD | 73.61 | 14.78 | 4.77 | 4.48 | 4.59 |
| | (Ours) PSLD ($\Gamma = 0.005$) | **44.36** | **6.99** | **3.77** | **3.92** | **4.17** |
| EM-US | CLD | 122.63 | 54.67 | 11.66 | 4.97 | 4.6 |
| | (Ours) PSLD ($\Gamma = 0.005$) | **99.05** | **44.06** | **8.9** | **4.42** | **4.37** |
| SSCS-QS | CLD | 44.83 | 10.7 | 4.82 | 4.73 | 4.74 |
| | (Ours) PSLD ($\Gamma = 0.005$) | **34.3** | **8.13** | **4.16** | **4.11** | **4.09** |
| SSCS-US | CLD | 105.16 | 45.54 | 6.75 | **4.06** | 4.18 |
| | (Ours) PSLD ($\Gamma = 0.005$) | **97.8** | **35.59** | **4.65** | 4.08 | **4.05** |

Table 14: Speed vs. Sample Quality comparison using the SDE setup for CelebA-64. FID computed for 10k samples. Values in **bold** indicate the best result for that column.

| Model | Size | NFE | FID $\downarrow$ | IS $\uparrow$ |
|---|---|---|---|---|
| PSLD ($\Gamma$=0.01) | 55M | 1000 | 2.34 | 9.57 |
| PSLD ($\Gamma$=0.02) | 55M | 1000 | 2.3 | 9.68 |
| PSLD ($\Gamma$=0.01, deep) | 97M | 1000 | 2.26 | 9.71 |
| PSLD ($\Gamma$=0.02, deep) | 97M | 1000 | **2.21** | **9.74** |

Table 15: CIFAR-10 sample quality (SDE). FID (lower is better) and IS (higher is better) were computed on 50k samples.

| Model | Size | NFE | FID $\downarrow$ | IS $\uparrow$ |
|---|---|---|---|---|
| PSLD ($\Gamma$=0.01) | 55M | 243 | 2.41 | 9.63 |
| PSLD ($\Gamma$=0.02) | 55M | 232 | 2.4 | 9.84 |
| PSLD ($\Gamma$=0.01, deep) | 97M | 246 | **2.10** | 9.79 |
| PSLD ($\Gamma$=0.02, deep) | 97M | 231 | 2.31 | 9.91 |
| PSLD ($\Gamma$=0.01, deep) | 97M | 159 | 2.13 | 9.76 |
| PSLD ($\Gamma$=0.02, deep) | 97M | 159 | 2.34 | **9.93** |

Table 16: CIFAR-10 sample quality (ODE). FID (lower is better) and IS (higher is better) were computed on 50k samples.

## D.3 Extended SOTA Results

**Extended Qualitative Results**: We provide qualitative samples from our SOTA CIFAR-10 models using the SDE and ODE setups in Figures 8 and 9 respectively. We provide some additional samples from the AFHQv2 dataset at the 128 x 128 resolution in Figure 10.

**Extended Quantitative Results**: Table 15 shows the FID and IS scores for all models using the SDE sampling setup. PSLD with $\Gamma = 0.02$ attains the best IS score of 9.74. When using the ODE sampling setup (Table 16), PSLD with $\Gamma = 0.02$ achieves the best IS score of 9.93.

## D.4 Conditional Synthesis using PSLD

**Class-Conditional Synthesis**: As discussed in Section 4.4, given class label information $\mathbf{y}$, an unconditional pre-trained score network $\boldsymbol{s}_\theta(\mathbf{z}_t, t)$ can be used for sampling from the class conditional distribution $p(\mathbf{z}_t|\mathbf{y})$ in PSLD. More specifically, we need to simulate the following reverse SDE:

$$d\mathbf{z}_t = \left[\boldsymbol{f}(\mathbf{z}_t) - \boldsymbol{G}(t)\boldsymbol{G}(t)^T \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{y})\right] dt + \boldsymbol{G}(t)d\mathbf{w}_t \tag{143}$$

The *conditional* score $\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{y})$ can be further decomposed as follows:

$$\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{y}) = \nabla_{\mathbf{z}_t} \log p(\mathbf{y}|\mathbf{z}_t) + \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \tag{144}$$

$$\approx \underbrace{\nabla_{\mathbf{z}_t} \log p(\mathbf{y}|\mathbf{z}_t)}_{\text{Classifier Gradient}} + \underbrace{\boldsymbol{s}_\theta(\mathbf{z}_t, t)}_{\text{Score}} \tag{145}$$

For practical scenarios, it is common to scale the contribution of the classifier gradient by a factor of $\lambda > 1$. Thus,

$$\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{y}) = \lambda \nabla_{\mathbf{z}_t} \log p(\mathbf{y}|\mathbf{z}_t) + \boldsymbol{s}_\theta(\mathbf{z}_t, t) \tag{146}$$

The above technique of approximating the conditional score $\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t|\mathbf{y})$ is called as *classifier-guidance* [4, 5]. The classifier $p(\mathbf{y}|\mathbf{z}_t)$ is trained by minimizing a time-dependent cross-entropy loss as follows:

$$\mathcal{L}_{\text{clf}}(\phi) = \mathbb{E}_{t\sim\mathcal{U}(0,1)}\mathbb{E}_{\mathbf{x}_0,\mathbf{y}\sim p_{\text{data}(\mathbf{x}_0,\mathbf{y})}}\mathbb{E}_{\mathbf{z}_t\sim p(\mathbf{z}_t|\mathbf{x}_0)}\left[-\sum_k \mathbb{1}(y = y_k)\log C_\phi^k(\mathbf{z}_t, t)\right] \tag{147}$$

where $C_\phi^k(\mathbf{z}_t, t)$ is a time-dependent classifier that takes as input a perturbed data point $\mathbf{z}_t$ and outputs class prediction probabilities. We perform class conditional synthesis for the CIFAR-10 (10 classes) and the AFHQ-v2 datasets. For the AFHQ-v2 dataset, we use the classes *Cats*, *Dogs*, and *Others* from the train split for classifier training (See Appendix C.6 for complete implementation details). We provide additional class conditional samples for CIFAR-10 in Figure 11 and for AFHQ-v2 in Figure 12.

**Image Inpainting**: Following [4], we can partition the input $\mathbf{x}_0$ into known ($\hat{\mathbf{x}}_0$) and unknown ($\bar{\mathbf{x}}_0$) components respectively. We can now define the diffusion for the unknown component in the augmented space as follows:

$$d\bar{\mathbf{z}}_t = \bar{\mathbf{f}}(\mathbf{z}_t)dt + \bar{\boldsymbol{G}}(t)d\mathbf{w}_t \tag{148}$$

where $\bar{\mathbf{f}}(\mathbf{z}_t) = \mathbf{f}(\bar{\mathbf{z}}_t)$ i.e. the drift applied to the missing components of $\mathbf{z}_t$. Similarly, $\bar{\boldsymbol{G}}(t)$ corresponds to the diffusion coefficient applied to the corresponding components of Brownian motion $d\mathbf{w}_t$. The corresponding reverse-SDE (conditioned on the observed signal $\hat{\mathbf{x}}_0$) can be specified as:

$$d\bar{\mathbf{z}}_t = \left[ \bar{\mathbf{f}}(\mathbf{z}_t) - \bar{\boldsymbol{G}}(t)\bar{\boldsymbol{G}}^T(t)\nabla_{\bar{\mathbf{z}}_t} \log p(\bar{\mathbf{z}}_t|\hat{\mathbf{x}}_0) \right] dt + \bar{\boldsymbol{G}}(t)d\mathbf{w}_t \tag{149}$$

Following the derivation in [4], it can be shown that:

$$\nabla_{\bar{\mathbf{z}}_t} \log p(\bar{\mathbf{z}}_t|\hat{\mathbf{x}}_0) \approx \nabla_{\bar{\mathbf{z}}_t} \log p(\bar{\mathbf{z}}_t|\hat{\mathbf{z}}_t) \tag{150}$$

$$= \nabla_{\bar{\mathbf{z}}_t} \log p([\bar{\mathbf{z}}_t; \hat{\mathbf{z}}_t]) \tag{151}$$

where $\hat{\mathbf{z}}_t \sim p(\hat{\mathbf{z}}_t|\hat{\mathbf{x}}_0)$ is a noisy augmented state sampled from the perturbation kernel given an observed signal $\hat{\mathbf{x}}_0$. We provide additional imputation results in Figure 13

**General Inverse Problems**: Similar to imputation, we can utilize PSLD for solving general inverse problems. Given a conditioning signal $\mathbf{y}$, we have,

$$\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t \mid \mathbf{y}) = \nabla_{\mathbf{z}_t} \log \int p_t(\mathbf{z}_t \mid \mathbf{y}_t, \mathbf{y})p(\mathbf{y}_t \mid \mathbf{y})d\mathbf{y}_t, \tag{152}$$

Further assuming that $p(\mathbf{y}_t \mid \mathbf{y})$ is tractable and $p_t(\mathbf{z}_t \mid \mathbf{y}_t, \mathbf{y}) \approx p_t(\mathbf{z}_t \mid \mathbf{y}_t)$, we have

$$\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t \mid \mathbf{y}) \approx \nabla_{\mathbf{z}_t} \log \int p_t(\mathbf{z}_t \mid \mathbf{y}_t)p(\mathbf{y}_t \mid \mathbf{y})d\mathbf{y}_t \tag{153}$$
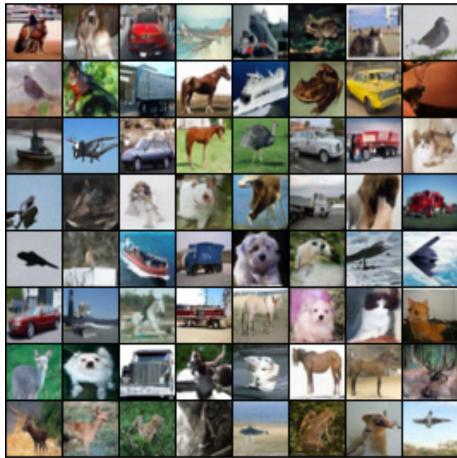
$$\approx \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t \mid \hat{\mathbf{y}}_t) \tag{154}$$

$$= \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) + \nabla_{\mathbf{z}_t} \log p_t(\hat{\mathbf{y}}_t \mid \mathbf{z}_t) \tag{155}$$

$$\approx \boldsymbol{s}_{\boldsymbol{\theta}^*}(\mathbf{z}_t, t) + \nabla_{\mathbf{z}_t} \log p_t(\hat{\mathbf{y}}_t \mid \mathbf{z}_t), \tag{156}$$
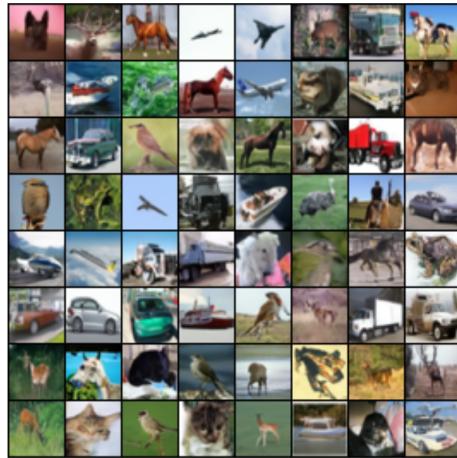
where $\hat{\mathbf{y}}_t \sim p(\mathbf{y}_t \mid \mathbf{y})$. Thus PSLD can be used for conditional synthesis like previous SGMs [4] while achieving better speed-vs-quality tradeoffs and better overall sample quality. Therefore, PSLD provides an attractive baseline for further developments in SGMs.
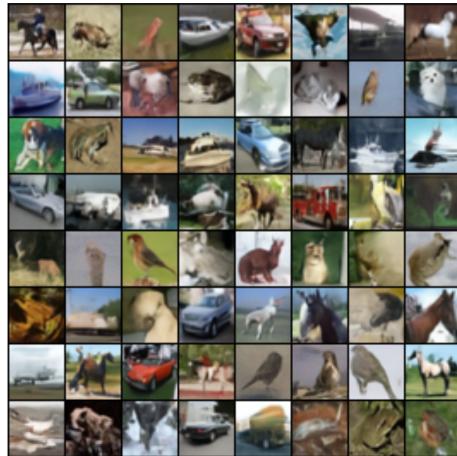
EM-US

EM-QS



$\Gamma = 0.25, \nu = 4.25$



$\Gamma = 0.25, \nu = 4.25$
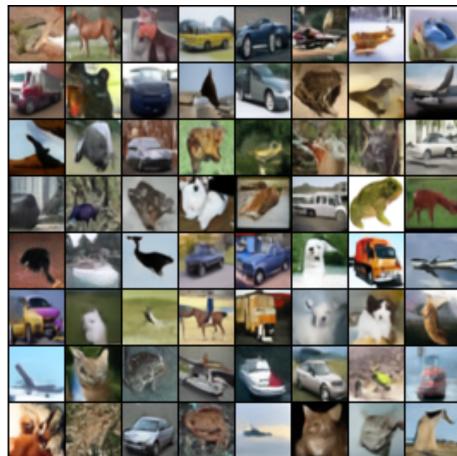


$\Gamma = 4.0, \nu = 0.0$



$\Gamma = 4.0, \nu = 0.0$



$\Gamma = 8.0, \nu = 4.0$



$\Gamma = 8.0, \nu = 4.0$

Figure 6: Qualitative illustration of the impact of $\Gamma$ on CIFAR-10 sample quality. Samples get progressively worse when increasing $\Gamma$. (Uncurated) Samples in the left and right columns were generated using EM-US and EM-QS samplers.
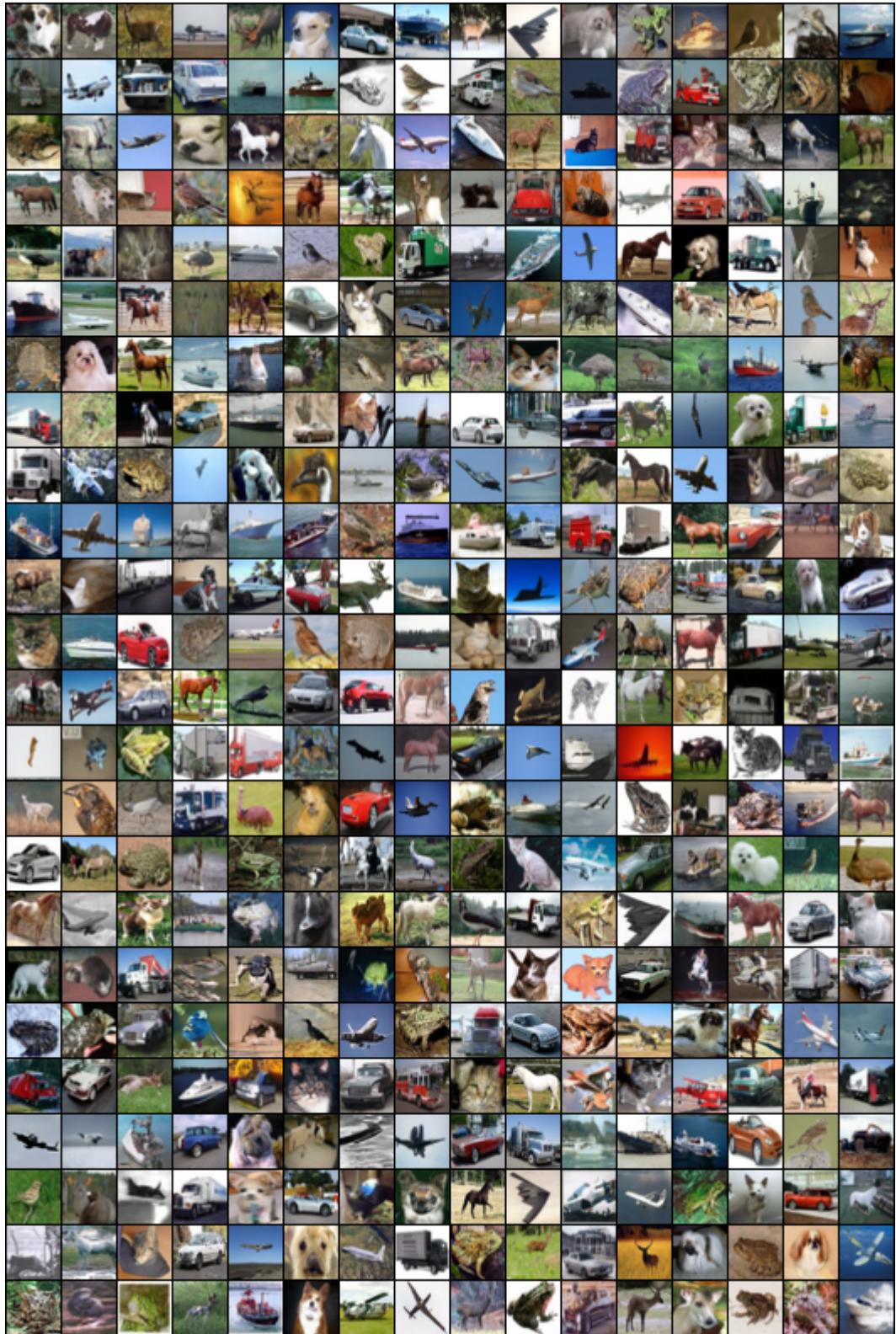
EM-US                                    EM-QS



$\Gamma = 0.005, \nu = 4.005$                $\Gamma = 0.005, \nu = 4.005$

$\Gamma = 0.02, \nu = 4.02$                  $\Gamma = 0.02, \nu = 4.02$

$\Gamma = 0.25, \nu = 4.25$                  $\Gamma = 0.25, \nu = 4.25$

Figure 7: Qualitative illustration of the impact of $\Gamma$ on CelebA-64 sample quality. Samples get progressively worse when increasing $\Gamma$. (Uncurated) Samples in the left and right columns were generated using EM-US and EM-QS samplers.

Figure 8: Uncurated samples from our SOTA PSLD ($\Gamma = 0.02, \nu = 4.02$) model using SDE sampling (FID=2.21, NFE=1000)
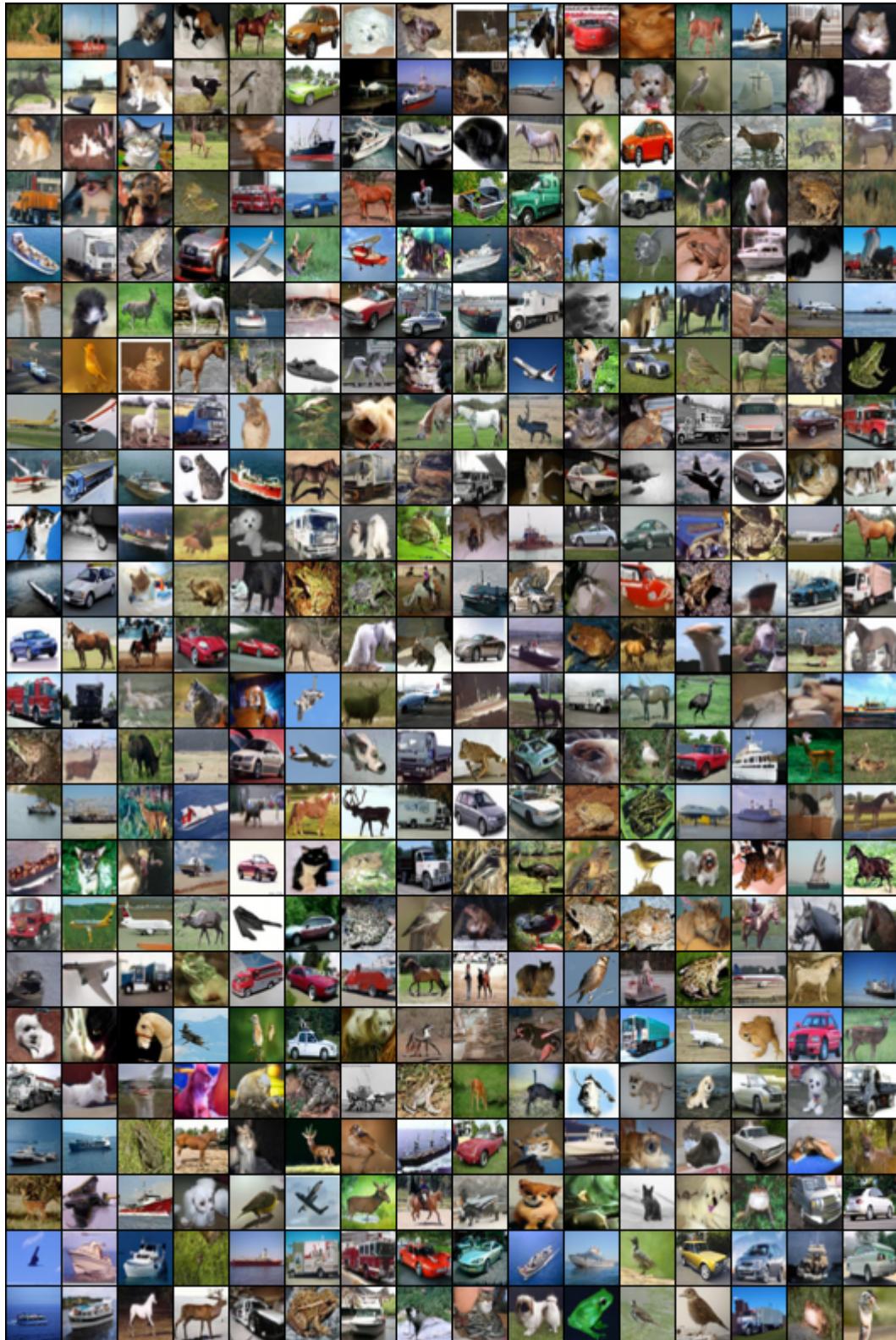
Figure 9: Uncurated samples from our SOTA PSLD ($\Gamma = 0.01, \nu = 4.01$) model using ODE sampling (FID=2.10, NFE=246)

Figure 10: Random unconditional AFHQv2 samples at 128x128 resolution from our PSLD ($\Gamma = 0.01$) model using the EM-QS sampler with N=1000.
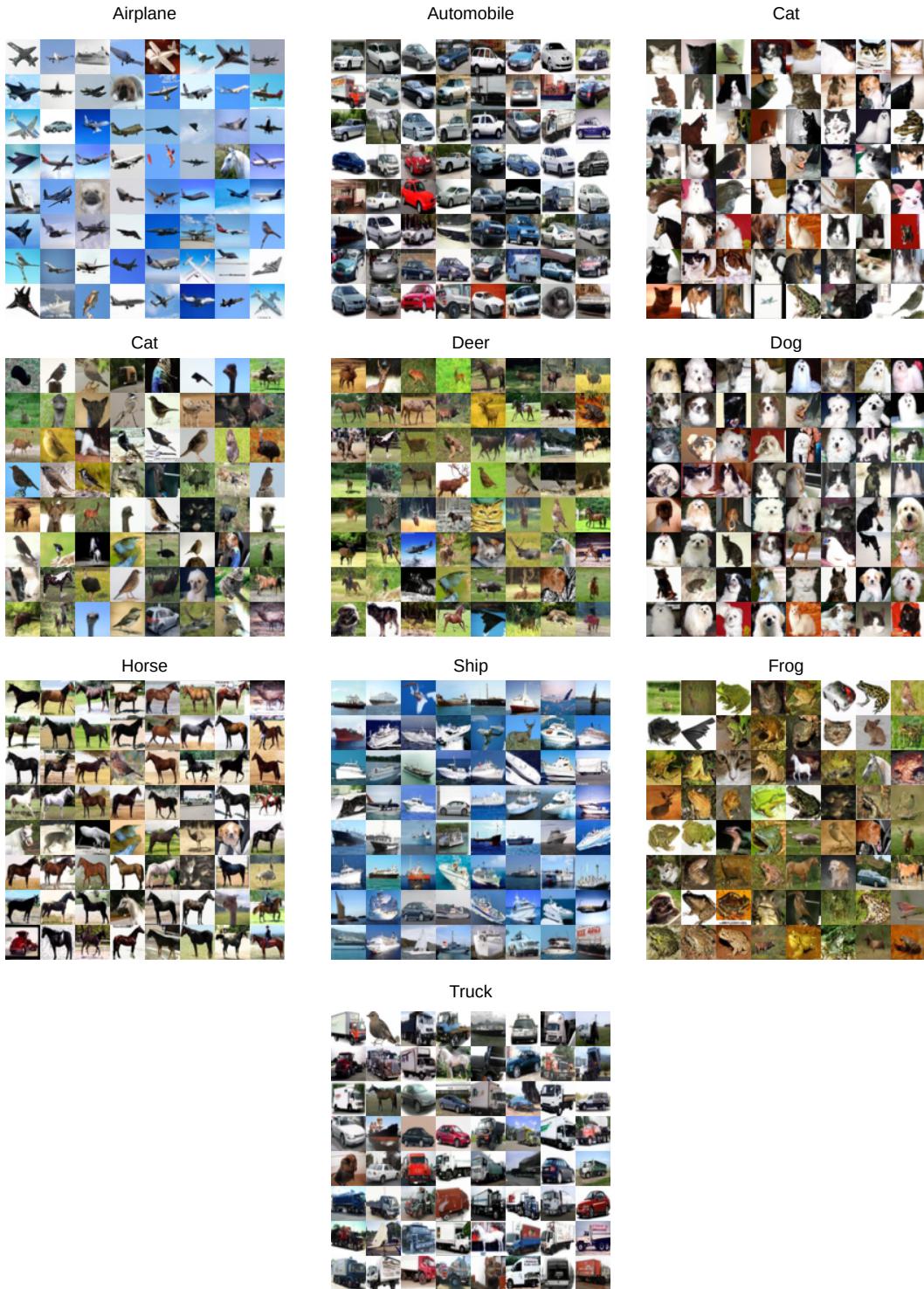
Airplane    Automobile    Cat

Cat    Deer    Dog

Horse    Ship    Frog

Truck

Figure 11: Randomly sampled class conditional results on the CIFAR-10 dataset using the EM-US sampler (N=1000). Guidance weight $\lambda = 5.0$. Using large guidance weight reduces diversity but improves sample quality.

Figure 12: Randomly sampled class conditional results on the AFHQv2 dataset using the EM-US sampler (N=1000). Guidance weight $\lambda = 10.0$. (Top to Bottom) Each of the three rows correspond to *Dog*, *Cats* and *Others*.

Figure 13: Additional imputation results on the AFHQv2 dataset (test split) using the EM-US sampler (N=1000). The Rightmost column indicates some failure cases.