
Diffusion Models With Learned Adaptive Noise

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Diffusion models have gained traction as powerful algorithms for synthesizing
2 high-quality images. Central to these algorithms is the diffusion process, which
3 maps data to noise according to equations inspired by thermodynamics and can
4 significantly impact performance. A widely held assumption is that the ELBO
5 objective of a diffusion model is invariant to the noise process (Kingma et al.,
6 2021). In this work, we dispel this assumption—we propose multivariate learned
7 adaptive noise (MULAN), a learned diffusion process that applies Gaussian noise
8 at different rates across an image. Our method consists of three components—a
9 multivariate noise schedule, instance-conditional diffusion, and auxiliary variables—
10 which ensure that the learning objective is no longer invariant to the choice of the
11 noise schedule as in previous works. Our work is grounded in Bayesian inference
12 and casts the learned diffusion process as an approximate variational posterior that
13 yields a tighter lower bound on marginal likelihood. Empirically, MULAN sets a
14 new state-of-the-art in density estimation on CIFAR-10 and ImageNet compared to
15 classical diffusion.

16 1 Introduction

17 A diffusion process q transforms an input datapoint denoted by \mathbf{x}_0 and sampled from a distribution
18 $q(\mathbf{x}_0)$ into a sequence of noisy data instances \mathbf{x}_t for $t \in [0, 1]$ by progressively adding Gaussian
19 noise of increasing magnitude. (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020). The
20 marginal distribution of each latent is defined by $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \alpha_t\mathbf{x}_0, \sigma_t\mathbf{I})$ where the diffusion
21 parameters $\alpha_t, \sigma_t \in \mathbb{R}^+$ implicitly define a noise schedule as a function of t , such that $\nu(t) = \alpha_t^2/\sigma_t^2$
22 is a monotonically decreasing function in t . Given any discretization of time into T timesteps of
23 width $1/T$, we define $t(i) = i/T$ and $s(i) = (i-1)/T$ and we use $\mathbf{x}_{0:1}$ to denote the subset of
24 variables associated with these timesteps; the forward process q can be shown to factorize into a
25 Markov chain $q(\mathbf{x}_{0:1}) = q(\mathbf{x}_0) \left(\prod_{i=1}^T q(\mathbf{x}_{t(i)}|\mathbf{x}_{s(i)}) \right)$.

26 The diffusion model p_θ is defined by a neural network (with parameters θ) used to denoise
27 the forward process q . Given a discretization of time into T steps, p factorizes as $p_\theta(\mathbf{x}_{0:1}) =$
28 $p_\theta(\mathbf{x}_1) \prod_{i=1}^T p_\theta(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)})$. We treat the \mathbf{x}_t for $t > 0$ as latent variables and fit p_θ by maximizing
29 the evidence lower bound (ELBO) on the marginal log-likelihood given by:

$$\log p_\theta(\mathbf{x}_0) = \text{ELBO}(p_\theta, q) + \text{D}_{\text{KL}}[q(\mathbf{x}_{t(1):t(T)}|\mathbf{x}_0)||p_\theta(\mathbf{x}_{t(1):t(T)}|\mathbf{x}_0)] \geq \text{ELBO}(p_\theta, q) \quad (1)$$

30 In most works, the noise schedule, as defined by $\nu(t)$, is either fixed or treated as a hyper-
31 parameter (Ho et al., 2020; Chen, 2023; Hoogeboom et al., 2023). Chen (2023); Hoogeboom
32 et al. (2023) show that the noise schedule can have a significant impact on sample quality. Kingma
33 et al. (2021) consider learning $\nu(t)$, but argue that the KL divergence terms in the ELBO are invariant
34 to the choice of function ν , except for the initial values $\nu(0), \nu(1)$, and they set these values to

35 hand-specified constants in their experiments. They only consider learning ν for the purpose of
 36 minimizing the variance of the gradient of the ELBO. In this work, we show that the ELBO is not
 37 invariant to more complex forward processes.

38 2 Diffusion Models With Multivariate Learned Adaptive Noise

39 **Why Learned Diffusion?** Perhaps the most direct motivation for our work comes from Bayesian
 40 inference. Notice that the gap between the evidence lower bound $\text{ELBO}(p, q)$ and the marginal log-
 41 likelihood (MLL) in Eq. 1 is precisely the KL divergence $D_{\text{KL}}[q(\mathbf{x}_{t(1):t(T)}|\mathbf{x}_0)||p_\theta(\mathbf{x}_{t(1):t(T)}|\mathbf{x}_0)]$
 42 between the diffusion process q over the latents \mathbf{x}_t and the true posterior of the diffusion model.
 43 The diffusion process takes the role of an approximate variational posterior in $\text{ELBO}(p, q)$. This
 44 observation suggests that the ELBO can be made tighter by choosing a diffusion processes q that
 45 is closer to the true posterior $p_\theta(\mathbf{x}_{t(1):t(T)}|\mathbf{x}_0)$; this in turn brings the learning objective of closer
 46 to $\log p(\mathbf{x})$, which is often the ideal objective that we wish to optimize. In fact, the key idea of
 47 variational inference is to optimize $\max_{q \in \mathcal{Q}} \text{ELBO}(p, q)$ over a family of approximate posteriors
 48 \mathcal{Q} to induce a tighter ELBO (Kingma & Welling, 2013). Most diffusion algorithms, however
 49 optimize $\max_{p \in \mathcal{P}} \text{ELBO}(p, q)$ within some family \mathcal{P} with a fixed q . Our work seeks to jointly
 50 optimize $\max_{p \in \mathcal{P}, q \in \mathcal{Q}} \text{ELBO}(p, q)$; we will show in our experiments that this improves the likelihood
 51 estimation.

52 2.1 A Forward Diffusion Process With Multivariate Adaptive Noise

53 This subsection focuses on defining \mathcal{Q} ; the next sections will show how to parameterize and train a
 54 reverse model $p \in \mathcal{P}$.

55 **Notation.** Given two vectors \mathbf{a} and \mathbf{b} , we use the notation $\mathbf{a}\mathbf{b}$ to represent the Hadamard product
 56 (element-wise multiplication). Additionally, we denote element-wise division of \mathbf{a} by \mathbf{b} as \mathbf{a}/\mathbf{b} . We
 57 denote the mapping $\text{diag}(\cdot)$ that takes a vector as input and produces a diagonal matrix as output.

58 2.1.1 Multivariate Gaussian Noise Schedule Conditioned on Context

59 Formally, our definition of a forward diffusion process with a multivariate noise schedule fol-
 60 lows previous work (Kingma et al., 2021; Hooeboom & Salimans, 2022) and defines q via
 61 the marginal for each latent noise variable \mathbf{x}_t for $t \in [0, 1]$, where the marginal is given by:
 62 $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\alpha}_t\mathbf{x}_0, \text{diag}(\boldsymbol{\sigma}_t^2))$, where $\mathbf{x}_t, \mathbf{x}_0 \in \mathbb{R}^d$, $\boldsymbol{\alpha}_t, \boldsymbol{\sigma}_t \in \mathbb{R}_+^d$ and d is the dimension-
 63 ality of the input data. For more details please refer Sec. B. In Sec. D.4, we argue that this class
 64 of diffusion process \mathcal{Q} induces an ELBO that is not invariant to $q \in \mathcal{Q}$. The ELBO consists of
 65 a line integral along the diffusion trajectory specified by $\nu(t)$. A line integrand is almost always
 66 path-dependent, unless its integral corresponds to a conservative force field, which is rarely the case
 67 for a diffusion process (Spinney & Ford, 2012). See Sec. D.4 for details.

68 Next, we extend the diffusion process to support context-adaptive noise. This enables injecting noise
 69 in a way that is dependent on the features of an image. Formally, we introduce a context variable
 70 $\mathbf{c} \in \mathbb{R}^m$ which encapsulates high-level information regarding \mathbf{x}_0 . Examples of \mathbf{c} could be a class
 71 label, a vector of attributes (e.g., features characterizing a human face), or even the input \mathbf{x}_0 itself. We
 72 define the marginal of the latent \mathbf{x}_t in the forward process as $q(\mathbf{x}_t|\mathbf{x}_0, \mathbf{c}) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\alpha}_t(\mathbf{c})\mathbf{x}_0, \boldsymbol{\sigma}_t^2(\mathbf{c}))$;
 73 the reverse process kernel can be similarly derived as Hooeboom & Salimans (2022):

$$q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0, \mathbf{c}) = \mathcal{N}\left(\boldsymbol{\mu}_q = \frac{\boldsymbol{\alpha}_{t|s}(\mathbf{c})\boldsymbol{\sigma}_s^2(\mathbf{c})}{\boldsymbol{\sigma}_t^2(\mathbf{c})}\mathbf{x}_t + \frac{\boldsymbol{\sigma}_{t|s}^2(\mathbf{c})\boldsymbol{\alpha}_s(\mathbf{c})}{\boldsymbol{\sigma}_t^2(\mathbf{c})}\mathbf{x}_0, \boldsymbol{\Sigma}_q = \text{diag}\left(\frac{\boldsymbol{\sigma}_s^2(\mathbf{c})\boldsymbol{\sigma}_{t|s}^2(\mathbf{c})}{\boldsymbol{\sigma}_t^2(\mathbf{c})}\right)\right) \quad (2)$$

74 where the diffusion parameters $\boldsymbol{\alpha}_t, \boldsymbol{\sigma}_t$ are now conditioned on \mathbf{c} via a neural network. Specifically,
 75 we parameterize the diffusion parameters $\boldsymbol{\alpha}_t(\mathbf{c}), \boldsymbol{\sigma}_t(\mathbf{c}), \nu(t, \mathbf{c})$ as $\boldsymbol{\alpha}_t^2(\mathbf{c}) = \text{sigmoid}(-\gamma_\phi(\mathbf{c}, t))$,
 76 $\boldsymbol{\sigma}_t^2(\mathbf{c}) = \text{sigmoid}(\gamma_\phi(\mathbf{c}, t))$, and $\nu(\mathbf{c}, t) = \exp(-\gamma_\phi(\mathbf{c}, t))$. Here, $\gamma_\phi(\mathbf{c}, t) : \mathbb{R}^m \times [0, 1] \rightarrow$
 77 $[\gamma_{\min}, \gamma_{\max}]^d$ is a neural network with the property that $\gamma_\phi(\mathbf{c}, t)$ is monotonic in t . Following Kingma
 78 et al. (2021); Zheng et al. (2023), we set $\gamma_{\min} = -13.30$, $\gamma_{\max} = 5.0$. We express $\gamma_\phi(\mathbf{c}, t)$ as a
 79 monotonic degree 5 polynomial in t . Details about the exact functional form of this polynomial and
 80 its implementation can be found in Suppl. D.2. More such parameterizations are discussed in C.

81 **2.2 Auxiliary-Variable Reverse Diffusion Processes**

82 We introduce a class of approximate reverse processes \mathcal{P} that match the structure of \mathcal{Q} and that are
 83 naturally suited to the joint optimization $\max_{p \in \mathcal{P}, q \in \mathcal{Q}} \text{ELBO}(p, q)$.

84 Formally, we define a diffusion model where the reverse diffusion process is conditioned on the
 85 context \mathbf{c} . Specifically, given any discretization of $t \in [0, 1]$ into T time steps as in Sec. 1, we
 86 introduce a context-conditional diffusion model $p_\theta(\mathbf{x}_{0:1}|\mathbf{c})$ that factorizes as the Markov chain

$$p_\theta(\mathbf{x}_{0:1}|\mathbf{c}) = p_\theta(\mathbf{x}_1|\mathbf{c}) \prod_{i=1}^T p_\theta(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}, \mathbf{c}). \quad (3)$$

87 Given that the true reverse process is a Gaussian as specified in Eq. 2, the ideal p_θ matches this
 88 parameterization (the proof mirrors that of regular diffusion models; Suppl. C), which yields

$$p_\theta(\mathbf{x}_s|\mathbf{c}, \mathbf{x}_t) = \mathcal{N}\left(\boldsymbol{\mu}_p = \frac{\boldsymbol{\alpha}_{t|s}(\mathbf{c})\boldsymbol{\sigma}_s^2(\mathbf{c})}{\boldsymbol{\sigma}_t^2(\mathbf{c})}\mathbf{x}_t + \frac{\boldsymbol{\sigma}_{t|s}^2(\mathbf{c})\boldsymbol{\alpha}_s(\mathbf{c})}{\boldsymbol{\sigma}_t^2(\mathbf{c})}\mathbf{x}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_p = \text{diag}\left(\frac{\boldsymbol{\sigma}_s^2(\mathbf{c})\boldsymbol{\sigma}_{t|s}^2(\mathbf{c})}{\boldsymbol{\sigma}_t^2(\mathbf{c})}\right)\right), \quad (4)$$

89 where $\mathbf{x}_\theta(\mathbf{x}_t, t)$, is a neural network that approximates \mathbf{x}_0 . Instead of parameterizing $\mathbf{x}_\theta(\mathbf{x}_t, t)$
 90 directly using a neural network, we consider 2 other parameterizations. One is the noise para-
 91 meterization (Ho et al., 2020) where $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{c}, t)$ is the denoising model which is parameter-
 92 ized as $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) = (\mathbf{x}_t - \boldsymbol{\alpha}_t(\mathbf{c})\mathbf{x}_\theta(\mathbf{x}_t, t, \mathbf{c}))/\boldsymbol{\sigma}_t(\mathbf{c})$; see Suppl. D.1.1 and the other is velocity
 93 parameterization (Salimans & Ho, 2022) where $\mathbf{v}_\theta(\mathbf{x}_t, \mathbf{c}, t)$ is a neural network that models
 94 $\mathbf{v}_\theta(\mathbf{x}_t, \mathbf{c}, t) = (\boldsymbol{\alpha}_t(\mathbf{c})\mathbf{x}_t - \mathbf{x}_\theta(\mathbf{x}_t, \mathbf{c}, t))/\boldsymbol{\sigma}_t(\mathbf{c})$; see Suppl. D.1.2.

95 **2.2.1 Conditioning Noise on an Auxiliary Latent Variable**

96 In Suppl. C.3, we highlight the challenges when \mathbf{c} is deterministic, and hence, propose an alternative
 97 strategy for learning conditional forward and reverse processes p, q that feature the same structure
 98 and hence support efficient noise parameterization. Our approach is based on the introduction of
 99 auxiliary variables (Wang et al., 2023), which lift the distribution p_θ into an augmented latent space.

100 Specifically, we define $\mathbf{z} \in \mathbb{R}^m$ as a low-dimensional auxiliary latent variable and define a lifted
 101 $p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{x}|\mathbf{z})p_\theta(\mathbf{z})$, where $p_\theta(\mathbf{x}|\mathbf{z})$ is the conditional diffusion model from Eq. 3 (with context
 102 \mathbf{c} set to \mathbf{z}) and $p_\theta(\mathbf{z})$ is a simple prior (e.g., unit Gaussian or fully factored Bernoulli). The latents \mathbf{z}
 103 can be interpreted as a high-level semantic representation of \mathbf{x} that conditions both the forward and
 104 the reverse processes. Unlike $\mathbf{x}_{0:1}$, the \mathbf{z} are not constrained to have a particular dimension and can
 105 be a low-dimensional vector of latent factors of variation. They can be continuous or discrete.

106 We form a learning objective for the lifted p_θ by applying the ELBO twice to obtain:

$$\log p_\theta(\mathbf{x}_0) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}_0)}[\log p_\theta(\mathbf{x}_0|\mathbf{z})] - \text{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}_0)||p_\theta(\mathbf{z})) \quad (5)$$

$$\geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}_0)}[\text{ELBO}(p_\theta(\mathbf{x}_{0:1}|\mathbf{z}), q_\phi(\mathbf{x}_{0:1}|\mathbf{z}))] - \text{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}_0)||p_\theta(\mathbf{z})), \quad (6)$$

107 where $\text{ELBO}(p_\theta(\mathbf{x}_{0:1}|\mathbf{z}), q_\phi(\mathbf{x}_{0:1}|\mathbf{z}))$ denotes the variational lower bound of a diffusion model
 108 (defined in Eq. 1) with a forward process $q_\phi(\mathbf{x}_{0:1}|\mathbf{z})$ (defined in Eq. 2 and Sec. 2.1.1) and an
 109 approximate reverse process $p_\theta(\mathbf{x}_{0:1}|\mathbf{z})$ (defined in Eq. 3), both conditioned on \mathbf{z} . The distribution
 110 $q_\phi(\mathbf{z}|\mathbf{x}_0)$ is an approximate posterior for \mathbf{z} parameterized by a neural network with parameters ϕ .

111 Crucially, note that in the learning objective (Eq. 6), the context, which in this case is \mathbf{z} , is available
 112 at training time in both the forward and reverse processes. At generation time, we can still obtain a
 113 valid context vector by sampling an auxiliary latent from $p_\theta(\mathbf{z})$. Thus, this approach addresses the
 114 aforementioned challenges and enables us to use the noise parameterization in Eq. 4.

115 **2.3 Variational Lower Bound**

116 Next, we derive a precise formula for the learning objective (6) of the auxiliary-variable diffusion
 117 model. Using the objective of a diffusion model in (1) we can write (6) as the sum of four terms:

$$\log p_\theta(\mathbf{x}_0) \geq \mathbb{E}_{q_\phi}[\mathcal{L}_{\text{recons}} + \mathcal{L}_{\text{diffusion}} + \mathcal{L}_{\text{prior}} + \mathcal{L}_{\text{latent}}], \quad (7)$$

118 The reconstruction loss, $\mathcal{L}_{\text{recons}}$, can be (stochastically and differentiably) estimated using stan-
 119 dard techniques; see (Kingma & Welling, 2013), $\mathcal{L}_{\text{prior}} = -\text{D}_{\text{KL}}[q_\phi(\mathbf{x}_1|\mathbf{x}_0, \mathbf{z})||p_\theta(\mathbf{x}_1)]$ is the

Table 1: Likelihood in bits per dimension (BPD) on the test set of CIFAR-10 and ImageNet. Results with “/” means they are not reported in the original papers. Model types are autoregressive (AR), normalizing flows (Flow), diffusion models (Diff). We only compare with results achieved without data augmentation.

Model	Type	CIFAR-10 (\downarrow)	ImageNet (\downarrow)
PixelCNN (Van den Oord et al., 2016)	AR	3.03	3.83
Image Transformer (Parmar et al., 2018)	AR	2.90	3.77
DDPM (Ho et al., 2020)	Diff	≤ 3.69	/
ScoreFlow (Song et al., 2021)	Diff	2.83	3.76
VDM (Kingma et al., 2021)	Diff	≤ 2.65	≤ 3.72
Flow Matching (Lipman et al., 2022)	Flow	2.99	/
Reflected Diffusion Models (Lou & Ermon, 2023)	Diff	2.68	3.74
MuLAN (Ours)	Diff	2.55 ± 10^{-3}	3.67 ± 10^{-3}

diffusion prior term, $\mathcal{L}_{\text{latent}} = -D_{\text{KL}}[q_{\phi}(\mathbf{z}|\mathbf{x}_0)||p_{\theta}(\mathbf{z})]$ is the latent prior term, and $\mathcal{L}_{\text{diffusion}} = -\frac{1}{2}\mathbb{E}_{t \sim [0,1]} [(\epsilon_t - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{z}, t))^{\top} \text{diag}(\nabla_t \gamma(\mathbf{z}, t)) (\epsilon_t - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{z}, t))]$ where $\nabla_t \gamma(\mathbf{z}, t) \in \mathbb{R}^d$ denotes the Jacobian of $\gamma(\mathbf{z}, t)$ with respect to the scalar t . We try two different kinds of priors for $p_{\theta}(\mathbf{z})$: discrete ($\mathbf{z} \in \{0, 1\}^m$) and continuous ($\mathbf{z} \in \mathbb{R}^m$). The exact expression for $\mathcal{L}_{\text{prior}}$ can be found in Suppl. D.3.

3 Experiments

This section reports experimental results on the CIFAR-10 (Krizhevsky et al., 2009) and ImageNet-32 (Van Den Oord et al., 2016) datasets. More details can be found in Sec. F.

Likelihood Estimation. In Table 1, we present likelihood estimation results for MuLAN and recent methods on CIFAR-10 and ImageNet-32 using the VLB-estimate; details in Sec. H.1. Trained for 10M steps on CIFAR-10 and 2M steps on ImageNet-32, this MuLAN version employs noise parameterization, akin to VDM (Kingma et al., 2021). Applied on VDM, MuLAN with a learned multivariate noising schedule conditioned on auxiliary latent variables significantly improves BPD over vanilla VDM. Additionally, using the ODE-based exact likelihood estimate, MuLAN outperforms existing methods in density estimation on both datasets, trained for 8M steps on CIFAR-10 and 2M steps on ImageNet-32. In inference, we extract the underlying probability flow ODE, similar to Zheng et al. (2023). While Zheng et al. (2023) used additional techniques, combining them with MuLAN could enhance its performance.

Ablations and Noise Schedule. In Fig. 3 we ablate each component of MuLAN. As $\gamma_{\phi}(\mathbf{z}, t)$ is multivariate, diverse noise schedules are expected for distinct input dimensions and $\mathbf{z} \sim p_{\theta}(\mathbf{z})$. In Fig. 2 with our CIFAR-10 model, we visualize the time-dependent variance of the noise schedule for different pixels, based on 128 samples $\mathbf{z} \sim p_{\theta}(\mathbf{z})$. Early schedule segments show increased variation, yet absolute variance is smaller than anticipated. Attempts to visualize noise schedules across diverse dataset images and areas (see Fig. 11) and with synthetic datasets exhibit no interpretable patterns, despite observable differences in likelihood estimation. We posit that alternative architectures and conditioning forms may unveil interpretable variations, a direction for future exploration.

4 Conclusion

In this study, we introduce MuLAN, a context-adaptive noise process that applies Gaussian noise at varying rates across input data. We present theoretical arguments challenging the prevailing notion that the likelihood of diffusion models remains independent of the noise schedules. We contend that this independence only holds true for univariate schedules, and in the case of multivariate schedules like MuLAN, different diffusion trajectories yield distinct likelihood estimates. Our evaluation of MuLAN spans multiple image datasets, where it outperforms state-of-the-art generative diffusion models. In general, MuLAN represents a promising avenue for advancing generative modeling.

References

- 154
155 Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary
156 differential equations. *Advances in neural information processing systems*, 31, 2018.
- 157 Ting Chen. On the importance of noise scheduling for diffusion models. *arXiv preprint*
158 *arXiv:2301.10972*, 2023.
- 159 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-
160 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,
161 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 162 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
163 *in Neural Information Processing Systems*, 34:8780–8794, 2021.
- 164 J.R. Dormand and P.J. Prince. A family of embedded runge-kutta formulae. *Journal of Computational*
165 *and Applied Mathematics*, 6(1):19–26, 1980. ISSN 0377-0427. doi: [https://doi.org/10.1016/](https://doi.org/10.1016/0771-050X(80)90013-3)
166 [0771-050X\(80\)90013-3](https://doi.org/10.1016/0771-050X(80)90013-3). URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0771050X80900133)
167 [0771050X80900133](https://www.sciencedirect.com/science/article/pii/S0771050X80900133).
- 168 Pavlos S. Efraimidis and Paul G. Spirakis. Weighted random sampling with a reservoir. *Infor-*
169 *mation Processing Letters*, 97(5):181–185, 2006. ISSN 0020-0190. doi: [https://doi.org/10.](https://doi.org/10.1016/j.ipl.2005.11.003)
170 [1016/j.ipl.2005.11.003](https://doi.org/10.1016/j.ipl.2005.11.003). URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S002001900500298X)
171 [S002001900500298X](https://www.sciencedirect.com/science/article/pii/S002001900500298X).
- 172 Will Grathwohl, Ricky TQ Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. Ffjord:
173 Free-form continuous dynamics for scalable reversible generative models. *arXiv preprint*
174 *arXiv:1810.01367*, 2018.
- 175 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
176 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 177 Emiel Hoogeboom and Tim Salimans. Blurring diffusion models. *arXiv preprint arXiv:2209.05557*,
178 2022.
- 179 Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for
180 high resolution images. *arXiv preprint arXiv:2301.11093*, 2023.
- 181 Michael F Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian
182 smoothing splines. *Communications in Statistics-Simulation and Computation*, 18(3):1059–1076,
183 1989.
- 184 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances*
185 *in neural information processing systems*, 34:21696–21707, 2021.
- 186 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
187 *arXiv:1312.6114*, 2013.
- 188 Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions.
189 *Advances in neural information processing systems*, 31, 2018.
- 190 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 191 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
192 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 193 Aaron Lou and Stefano Ermon. Reflected diffusion models. In *International Conference on Machine*
194 *Learning*, pp. 22675–22701. PMLR, 2023.
- 195 Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models, 2021.
- 196 Mathias Niepert, Pasquale Minervini, and Luca Franceschi. Implicit MLE: backpropagating through
197 discrete exponential family distributions. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N.
198 Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 2021, NeurIPS*
199 *2021, December 6-14, 2021, virtual*, pp. 14567–14579, 2021. URL [https://proceedings.](https://proceedings.neurips.cc/paper/2021/hash/7a430339c10c642c4b2251756fd1b484-Abstract.html)
200 [neurips.](https://proceedings.neurips.cc/paper/2021/hash/7a430339c10c642c4b2251756fd1b484-Abstract.html)
201 [cc/paper/2021/hash/7a430339c10c642c4b2251756fd1b484-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/7a430339c10c642c4b2251756fd1b484-Abstract.html).

- 202 Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and
203 Dustin Tran. Image transformer. In *International conference on machine learning*, pp. 4055–4064.
204 PMLR, 2018.
- 205 Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual
206 reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial*
207 *intelligence*, volume 32, 2018.
- 208 Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn. Diffu-
209 sion autoencoders: Toward a meaningful and decodable representation. In *IEEE Conference on*
210 *Computer Vision and Pattern Recognition (CVPR)*, 2022.
- 211 Subham Sekhar Sahoo, Anselm Paulus, Marin Vlastelica, Vít Musil, Volodymyr Kuleshov, and
212 Georg Martius. Backpropagation through combinatorial algorithms: Identity with projection
213 works. In *The Eleventh International Conference on Learning Representations*, 2023. URL
214 <https://openreview.net/forum?id=JZMR727029>.
- 215 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv*
216 *preprint arXiv:2202.00512*, 2022.
- 217 Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the
218 pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv preprint*
219 *arXiv:1701.05517*, 2017.
- 220 John Skilling. The eigenvalues of mega-dimensional matrices. 1989. URL [https://api.](https://api.semanticscholar.org/CorpusID:117844915)
221 [semanticscholar.org/CorpusID:117844915](https://api.semanticscholar.org/CorpusID:117844915).
- 222 Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
223 learning using nonequilibrium thermodynamics, 2015.
- 224 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
225 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
226 *arXiv:2011.13456*, 2020.
- 227 Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of
228 score-based diffusion models. *Advances in neural information processing systems*, 34:1415–1428,
229 2021.
- 230 Richard E. Spinney and Ian J. Ford. Fluctuation relations: a pedagogical overview, 2012.
- 231 Benigno Uribe, Iain Murray, and Hugo Larochelle. Rnade: The real-valued neural autoregressive
232 density-estimator. *Advances in Neural Information Processing Systems*, 26, 2013.
- 233 Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional
234 image generation with pixelcnn decoders. *Advances in neural information processing systems*, 29,
235 2016.
- 236 Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks.
237 In *International conference on machine learning*, pp. 1747–1756. PMLR, 2016.
- 238 Yingheng Wang, Yair Schiff, Aaron Gokaslan, Weishen Pan, Fei Wang, Christopher De Sa, and
239 Volodymyr Kuleshov. Infodiffusion: Representation learning using information maximizing
240 diffusion models. In *International Conference on Machine Learning*, pp. xxxx–xxxx. PMLR,
241 2023.
- 242 Sang Michael Xie and Stefano Ermon. Reparameterizable subset sampling via continuous relaxations.
243 *arXiv preprint arXiv:1901.10517*, 2019.
- 244 Kaiwen Zheng, Cheng Lu, Jianfei Chen, and Jun Zhu. Improved techniques for maximum likelihood
245 estimation for diffusion odes. *arXiv preprint arXiv:2305.03935*, 2023.

246 A Standard Diffusion Models

247 We have a Gaussian diffusion process that begins with the data \mathbf{x}_0 , and defines a sequence of
 248 increasingly noisy versions of \mathbf{x}_0 which we call the latent variables \mathbf{x}_t , where t runs from $t = 0$ (least
 249 noisy) to $t = 1$ (most noisy). Given, T , we discretize time uniformly into T timesteps each with a
 250 width $1/T$. We define $t(i) = i/T$ and $s(i) = (i - 1)/T$.

251 A.1 Forward Process

$$q(\mathbf{x}_t|\mathbf{x}_s) = \mathcal{N}(\alpha_{t|s}\mathbf{x}_s, \sigma_{t|s}^2\mathbf{I}_n) \quad (8)$$

252 where

$$\alpha_{t|s} = \frac{\alpha_t}{\alpha_s} \quad (9)$$

$$\sigma_{t|s}^2 = \sigma_t^2 - \frac{\alpha_{t|s}^2}{\sigma_s^2} \quad (10)$$

253 A.2 Reverse Process

254 [Kingma et al. \(2021\)](#) show that the distribution $q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0)$ is also gaussian,

$$q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\boldsymbol{\mu}_q = \frac{\alpha_{t|s}\sigma_s^2}{\sigma_t^2}\mathbf{x}_t + \frac{\sigma_{t|s}^2\alpha_s}{\sigma_t^2}\mathbf{x}_0, \boldsymbol{\Sigma}_q = \frac{\sigma_s^2\sigma_{t|s}^2}{\sigma_t^2}\mathbf{I}_n\right) \quad (11)$$

255 Since during the reverse process, we don't have access to \mathbf{x}_0 , we approximate it using a neural
 256 network $\mathbf{x}_\theta(\mathbf{x}_t, t)$ with parameters θ . Thus,

$$p_\theta(\mathbf{x}_s|\mathbf{x}_t) = \mathcal{N}\left(\boldsymbol{\mu}_p = \frac{\alpha_{t|s}\sigma_s^2}{\sigma_t^2}\mathbf{x}_t + \frac{\sigma_{t|s}^2\alpha_s}{\sigma_t^2}\mathbf{x}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_p = \frac{\sigma_s^2\sigma_{t|s}^2}{\sigma_t^2}\mathbf{I}_n\right) \quad (12)$$

257 A.3 Variational Lower Bound

258 This corruption process q is the following markov-chain as $q(\mathbf{x}_{0:1}) = q(\mathbf{x}_0) \left(\prod_{i=1}^T q(\mathbf{x}_{t(i)}|\mathbf{x}_{s(i)})\right)$.
 259 In the reverse Rrocess, or the denoising process, p_θ , a neural network (with parameters θ)
 260 is used to denoise the noising process q . The reverse Rrocess factorizes as: $p_\theta(\mathbf{x}_{0:1}) =$
 261 $p_\theta(\mathbf{x}_1) \prod_{i=1}^T p_\theta(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)})$. Let $\mathbf{x}_\theta(\mathbf{x}_t, t)$ be the reconstructed input by a neural network from
 262 \mathbf{x}_t . Similar to [Sohl-Dickstein et al. \(2015\)](#); [Kingma et al. \(2021\)](#) we decompose the negative lower
 263 bound (VLB) as:

$$\begin{aligned} -\log p_\theta(\mathbf{x}_0) &\leq \mathbb{E}_{q_\phi} \left[-\log \frac{p_\theta(\mathbf{x}_{t(0):t(T)})}{q_\phi(\mathbf{x}_{t(1):t(T)}|\mathbf{x}_0)} \right] \\ &= \mathbb{E}_{\mathbf{x}_{t(1)} \sim q(\mathbf{x}_{t(1)}|\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0|\mathbf{x}_{t(1)})] \\ &\quad + \sum_{i=2}^T \mathbb{E}_{\mathbf{x}_{t(i)}|\mathbf{x}_0} \text{D}_{\text{KL}}[p_\theta(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}) \| q_\phi(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}, \mathbf{x}_0)] \\ &\quad + \text{D}_{\text{KL}}[p_\theta(\mathbf{x}_1) \| q_\phi(\mathbf{x}_1|\mathbf{x}_0)] \\ &= \underbrace{\mathbb{E}_{\mathbf{x}_{t(1)} \sim q(\mathbf{x}_{t(1)}|\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0|\mathbf{x}_{t(1)})]}_{\mathcal{L}_{\text{recons}}} \\ &\quad + \underbrace{\frac{T}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n), i \sim \mathcal{U}\{2, T\}} \text{D}_{\text{KL}}[p_\theta(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}) \| q_\phi(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}, \mathbf{x}_0)]}_{\mathcal{L}_{\text{diffusion}}} \\ &\quad + \underbrace{\text{D}_{\text{KL}}[p_\theta(\mathbf{x}_1) \| q_\phi(\mathbf{x}_1|\mathbf{x}_0)]}_{\mathcal{L}_{\text{prior}}} \end{aligned} \quad (13)$$

264 The prior loss, $\mathcal{L}_{\text{prior}}$, and reconstruction loss, $\mathcal{L}_{\text{recons}}$, can be (stochastically and differentiably)
 265 estimated using standard techniques; see [Kingma & Welling \(2013\)](#). The diffusion loss, $\mathcal{L}_{\text{diffusion}}$,
 266 varies with the formulation of the noise schedule. We provide an exact formulation for it in the
 267 subsequent sections.

268 A.4 Diffusion Loss

269 For brevity, we use the notation s for $s(i)$ and t for $t(i)$. From Eq. 25 and Eq. 26 we get the following
 270 expression for $q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0)$:

$$\begin{aligned} & \text{D}_{\text{KL}}(q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_s|\mathbf{x}_t)) \\ &= \frac{1}{2} \left((\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^\top \boldsymbol{\Sigma}_{\theta}^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) + \text{tr}(\boldsymbol{\Sigma}_q \boldsymbol{\Sigma}_p^{-1} - \mathbf{I}_n) - \log \frac{|\boldsymbol{\Sigma}_q|}{|\boldsymbol{\Sigma}_p|} \right) \\ &= \frac{1}{2} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^\top \boldsymbol{\Sigma}_{\theta}^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) \end{aligned}$$

Substituting $\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q, \boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p$ from equation 12 and equation 11; for the exact derivation see [Kingma et al. \(2021\)](#)

$$= \frac{1}{2} (\nu(s) - \nu(t)) \|\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)\|_2^2 \quad (14)$$

271 Thus $\mathcal{L}_{\text{diffusion}}$ is given by

$$\begin{aligned} & \mathcal{L}_{\text{diffusion}} \\ &= \lim_{T \rightarrow \infty} \frac{T}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n), i \sim U\{2, T\}} \text{D}_{\text{KL}}[p_{\theta}(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}) \| q_{\phi}(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}, \mathbf{x}_0)] \\ &= \lim_{T \rightarrow \infty} \frac{1}{2} \sum_{i=2}^T \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} (\nu(s) - \nu(t)) \|\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)\|_2^2 \\ &= \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} \left[\lim_{T \rightarrow \infty} \sum_{i=2}^T (\nu(s) - \nu(t)) \|\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)\|_2^2 \right] \\ &= \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} \left[\lim_{T \rightarrow \infty} \sum_{i=2}^T T (\nu(s) - \nu(t)) \|\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)\|_2^2 \frac{1}{T} \right] \end{aligned}$$

Substituting $\lim_{T \rightarrow \infty} T(\nu(s) - \nu(t)) = \frac{d}{dt} \nu(t) \equiv \nu'(t)$; see [Kingma et al. \(2021\)](#)

$$= \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} \left[\int_0^1 \nu'(t) \|\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)\|_2^2 dt \right] \quad (15)$$

In practice instead of computing the integral is computed by MC sampling.

$$= -\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n), t \sim U[0, 1]} [\nu'(t) \|\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)\|_2^2] \quad (16)$$

272 B Multivariate noise schedule

273 Our proposed forward diffusion process progressively induces varying amounts of Gaussian noise
 274 across different areas of the image. We introduce two new components relative to previous work:
 275 multivariate noise scheduling and context-adaptive noise.

276 Intuitively, a multivariate noise schedule injects noise at different rates for each pixel of an input
 277 image. This enables adapting the diffusion process to spatial variations within the image. We will
 278 also see that this change is sufficient to make the ELBO no longer invariant in q .

279 Formally, our definition of a forward diffusion process with a multivariate noise schedule follows
 280 previous work ([Kingma et al., 2021](#); [Hooeboom & Salimans, 2022](#)) and defines q via the marginal
 281 for each latent noise variable \mathbf{x}_t for $t \in [0, 1]$, where the marginal is given by:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\alpha}_t \mathbf{x}_0, \text{diag}(\boldsymbol{\sigma}_t^2)), \quad (17)$$

282 where $\mathbf{x}_t, \mathbf{x}_0 \in \mathbb{R}^d$, $\alpha_t, \sigma_t \in \mathbb{R}_+^d$ and d is the dimensionality of the input data. The α_t, σ_t denote
 283 varying amounts of signal and associated with each component (i.e., each pixel) of \mathbf{x}_0 as a function
 284 of time $t(i)$. Similarly to Kingma et al. (2021), we may define the multivariate signal-to-noise ratio
 285 as $\nu(t) = \alpha_t^2 / \sigma_t^2$ and we choose α_t, σ_t such that $\nu(t)$ is monotonically decreasing in t along all
 286 dimensions and is differentiable in $t \in [0, 1]$. Let $\alpha_{t|s} = \alpha_t / \alpha_s$ and $\sigma_{t|s}^2 = \sigma_t^2 - \alpha_{t|s}^2 / \sigma_s^2$ with
 287 all operations applied elementwise. In Hoogeboom & Salimans (2022), show that these marginals
 288 induce transition kernels of the true reverse process between steps $s < t$ that are given by:

$$q(\mathbf{x}_s | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_s; \boldsymbol{\mu}_q = \frac{\alpha_{t|s} \sigma_s^2}{\sigma_t^2} \mathbf{x}_t + \frac{\sigma_{t|s}^2 \alpha_s}{\sigma_t^2} \mathbf{x}_0, \boldsymbol{\Sigma}_q = \text{diag} \left(\frac{\sigma_s^2 \sigma_{t|s}^2}{\sigma_t^2} \right) \right) \quad (18)$$

289 In Sec. D.4, we argue that this class of diffusion process \mathcal{Q} induces an ELBO that is not invariant to
 290 $q \in \mathcal{Q}$. The ELBO consists of a line integral along the diffusion trajectory specified by $\nu(t)$. A line
 291 integrand is almost always path-dependent, unless its integral corresponds to a conservative force
 292 field, which is rarely the case for a diffusion process (Spinney & Ford, 2012). See Sec. D.4 for
 293 details.

294 For a multivariate noise schedule we have $\alpha_t, \sigma_t \in \mathbb{R}^{n \times n}$ where $t \in [0, 1]$. α_t, σ_t are diagonal
 295 matrices. The timesteps s, t satisfy $0 \leq s < t \leq 1$. Furthermore, we use the following notations
 296 where arithmetic division represents element wise division between 2 diagonal matrices:

$$\alpha_{t|s} = \frac{\alpha_t}{\alpha_s} \quad (19)$$

$$\sigma_{t|s}^2 = \sigma_t^2 - \frac{\alpha_{t|s}^2}{\sigma_s^2} \quad (20)$$

297 B.1 Forward Process

$$q(\mathbf{x}_t | \mathbf{x}_s) = \mathcal{N} \left(\alpha_{t|s} \mathbf{x}_s, \sigma_{t|s}^2 \right) \quad (21)$$

298 **Change of variables.** We can write \mathbf{x}_t explicitly in terms of the signal-to-noise ratio, $\nu(t)$, and
 299 input \mathbf{x}_0 in the following manner:

$$\nu_t = \frac{\alpha_t^2}{\sigma_t^2}$$

We know $\alpha_t^2 = 1 - \sigma_t^2$ for Variance Preserving process; see Sec. 1.

$$\implies \frac{1 - \sigma_t^2}{\sigma_t^2} = \nu_t$$

$$\implies \sigma_t^2 = \frac{1}{1 + \nu_t} \quad \text{and} \quad \alpha_t^2 = \frac{\nu_t}{1 + \nu_t} \quad (22)$$

300

$$\nu_t = \frac{\alpha_t^2}{\sigma_t^2}$$

We know $\alpha_t^2 = 1 - \sigma_t^2$ for Variance Preserving process; see Sec. 1.

$$\implies \frac{1 - \sigma_t^2}{\sigma_t^2} = \nu_t$$

$$\implies \sigma_t^2 = \frac{1}{1 + \nu_t} \quad \text{and} \quad \alpha_t^2 = \frac{\nu_t}{1 + \nu_t} \quad (23)$$

301 Thus, we write \mathbf{x}_t in terms of the signal-to-noise ratio in the following manner:

$$\begin{aligned} \mathbf{x}_{\nu(t)} &= \alpha_t \mathbf{x}_0 + \sigma_t \boldsymbol{\epsilon}_t; \quad \boldsymbol{\epsilon}_t \sim \mathcal{N}(0, \mathbf{I}_n) \\ &= \frac{\sqrt{\nu(t)}}{\sqrt{1 + \nu(t)}} \mathbf{x}_0 + \frac{1}{\sqrt{1 + \nu(t)}} \boldsymbol{\epsilon}_t \end{aligned} \quad \text{Using Eq. 22} \quad (24)$$

302 **B.2 Reverse Process**

303 The distribution of \mathbf{x}_t given \mathbf{x}_s is given by:

$$q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\boldsymbol{\mu}_q = \frac{\boldsymbol{\alpha}_{t|s}\sigma_s^2}{\sigma_t^2}\mathbf{x}_t + \frac{\sigma_{t|s}^2\boldsymbol{\alpha}_s}{\sigma_t^2}\mathbf{x}_0, \boldsymbol{\Sigma}_q = \text{diag}\left(\frac{\sigma_s^2\sigma_{t|s}^2}{\sigma_t^2}\right)\right) \quad (25)$$

304 Let $\mathbf{x}_\theta(\mathbf{x}_t, t)$ be the neural network approximation for \mathbf{x}_0 . Then we get the following reverse process:

$$p_\theta(\mathbf{x}_s|\mathbf{x}_t) = \mathcal{N}\left(\boldsymbol{\mu}_p = \frac{\boldsymbol{\alpha}_{t|s}\sigma_s^2}{\sigma_t^2}\mathbf{x}_t + \frac{\sigma_{t|s}^2\boldsymbol{\alpha}_s}{\sigma_t^2}\mathbf{x}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_p = \text{diag}\left(\frac{\sigma_s^2\sigma_{t|s}^2}{\sigma_t^2}\right)\right) \quad (26)$$

305 **B.3 Diffusion Loss**

306 For brevity we use the notation s for $s(i)$ and t for $t(i)$. From Eq. 25 and Eq. 26 we get the following
307 expression for $q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0)$:

$$\begin{aligned} & D_{\text{KL}}(q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_s|\mathbf{x}_t)) \\ &= \frac{1}{2} \left((\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^\top \boldsymbol{\Sigma}_\theta^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) + \text{tr}(\boldsymbol{\Sigma}_q \boldsymbol{\Sigma}_p^{-1} - \mathbf{I}_n) - \log \frac{|\boldsymbol{\Sigma}_q|}{|\boldsymbol{\Sigma}_p|} \right) \\ &= \frac{1}{2} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^\top \boldsymbol{\Sigma}_\theta^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) \end{aligned}$$

Substituting $\boldsymbol{\mu}_q, \boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p$ from equation 26 and equation 25.

$$\begin{aligned} &= \frac{1}{2} \left(\frac{\sigma_{t|s}^2 \boldsymbol{\alpha}_s}{\sigma_t^2} \mathbf{x}_0 - \frac{\sigma_{t|s}^2 \boldsymbol{\alpha}_s}{\sigma_t^2} \mathbf{x}_\theta(\mathbf{x}_t, t) \right)^\top \text{diag}\left(\frac{\sigma_s^2 \sigma_{t|s}^2}{\sigma_t^2}\right)^{-1} \left(\frac{\sigma_{t|s}^2 \boldsymbol{\alpha}_s}{\sigma_t^2} \mathbf{x}_0 - \frac{\sigma_{t|s}^2 \boldsymbol{\alpha}_s}{\sigma_t^2} \mathbf{x}_\theta(\mathbf{x}_t, t) \right) \\ &= \frac{1}{2} (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t))^\top \text{diag}\left(\frac{\sigma_{t|s}^2 \boldsymbol{\alpha}_s}{\sigma_t^2}\right)^\top \text{diag}\left(\frac{\sigma_s^2 \sigma_{t|s}^2}{\sigma_t^2}\right)^{-1} \text{diag}\left(\frac{\sigma_{t|s}^2 \boldsymbol{\alpha}_s}{\sigma_t^2}\right) (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t)) \\ &= \frac{1}{2} (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t))^\top \text{diag}\left(\frac{\sigma_{t|s}^2 \boldsymbol{\alpha}_s}{\sigma_t^2} \odot \frac{\sigma_t^2}{\sigma_s^2 \sigma_{t|s}^2} \odot \frac{\sigma_{t|s}^2 \boldsymbol{\alpha}_s}{\sigma_t^2}\right) (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t)) \\ &= \frac{1}{2} (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t))^\top \text{diag}\left(\frac{\sigma_{t|s}^2 \boldsymbol{\alpha}_s^2}{\sigma_t^2 \sigma_s^2}\right) (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t)) \end{aligned}$$

Simplifying the expression using eq. 19 and eq. 20 we get,

$$= \frac{1}{2} (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t))^\top \text{diag}\left(\frac{\boldsymbol{\alpha}_s^2}{\sigma_s^2} - \frac{\boldsymbol{\alpha}_t^2}{\sigma_t^2}\right) (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t))$$

Using the relation $\boldsymbol{\nu}(t) = \boldsymbol{\alpha}_t^2 / \sigma_t^2$ we get,

$$= \frac{1}{2} (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t))^\top \text{diag}(\boldsymbol{\nu}(s) - \boldsymbol{\nu}(t)) (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, t)) \quad (27)$$

308 Like Kingma et al. (2021) we train the model in the continuous domain with $T \rightarrow \infty$.

$$\begin{aligned}
& \mathcal{L}_{\text{diffusion}} \\
&= \lim_{T \rightarrow \infty} \frac{1}{2} \sum_{i=2}^T \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} \text{D}_{\text{KL}}(q(\mathbf{x}_{s(i)} | \mathbf{x}_{t(i)}, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{s(i)} | \mathbf{x}_{t(i)})) \\
&= \lim_{T \rightarrow \infty} \frac{1}{2} \sum_{i=2}^T \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_{t(i)}, t(i)))^{\top} \text{diag}(\boldsymbol{\nu}_{s(i)} - \boldsymbol{\nu}_{t(i)}) (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_{t(i)}, t)) \\
&= \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} \left[\lim_{T \rightarrow \infty} \sum_{i=2}^T (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_{t(i)}, t(i)))^{\top} \text{diag}(\boldsymbol{\nu}_{s(i)} - \boldsymbol{\nu}_{t(i)}) (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_{t(i)}, t)) \right] \\
&= \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} \left[\lim_{T \rightarrow \infty} \sum_{i=2}^T T (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_{t(i)}, t(i)))^{\top} \text{diag}(\boldsymbol{\nu}_{s(i)} - \boldsymbol{\nu}_{t(i)}) (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_{t(i)}, t)) \frac{1}{T} \right] \\
&\text{Let } \lim_{T \rightarrow \infty} T(\boldsymbol{\nu}_{s(i)} - \boldsymbol{\nu}_{t(i)}) = \frac{d}{dt} \boldsymbol{\nu}(t) \text{ denote the scalar derivative of the vector } \boldsymbol{\nu}(t) \text{ w.r.t } t \\
&= \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} \left[\int_0^1 (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t))^{\top} \text{diag} \left(\frac{d}{dt} \boldsymbol{\nu}(t) \right) (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)) dt \right] \tag{28}
\end{aligned}$$

In practice instead of computing the integral is computed by MC sampling.

$$= -\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n), t \sim U[0, 1]} \left[(\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t))^{\top} \text{diag} \left(\frac{d}{dt} \boldsymbol{\nu}(t) \right) (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)) \right] \tag{29}$$

309 B.4 Vectorized Representation of the diffusion loss

310 Let $\boldsymbol{\nu}(t)$ be the vectorized representation of the diagonal entries of the matrix $\boldsymbol{\nu}(t)$. We can rewrite
311 the integral in eq. 28 in the following vectorized form where \odot denotes element wise multiplication
312 and $\langle \cdot, \cdot \rangle$ denotes dot product between 2 vectors.

$$\begin{aligned}
& \mathcal{L}_{\text{diffusion}} \\
&= -\frac{1}{2} \int_0^1 (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t))^{\top} \text{diag} \left(\frac{d}{dt} \boldsymbol{\nu}(t) \right) (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)) dt \\
&= -\frac{1}{2} \int_0^1 \left\langle (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)) \odot (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)), \frac{d}{dt} \boldsymbol{\nu}(t) \right\rangle dt \\
&\text{Using change of variables as mentioned in Sec. 2.1 we have} \\
&= -\frac{1}{2} \int_0^1 \left\langle (\mathbf{x}_0 - \tilde{\mathbf{x}}_{\theta}(\mathbf{x}_{\boldsymbol{\nu}(t)}, \boldsymbol{\nu}(t))) \odot (\mathbf{x}_0 - \tilde{\mathbf{x}}_{\theta}(\mathbf{x}_{\boldsymbol{\nu}(t)}, \boldsymbol{\nu}(t))), \frac{d}{dt} \boldsymbol{\nu}(t) \right\rangle dt \\
&\text{Let } \mathbf{f}_{\theta}(\mathbf{x}_0, \boldsymbol{\nu}(t)) = (\mathbf{x}_0 - \tilde{\mathbf{x}}_{\theta}(\mathbf{x}_{\boldsymbol{\nu}(t)}, \boldsymbol{\nu}(t))) \odot (\mathbf{x}_0 - \tilde{\mathbf{x}}_{\theta}(\mathbf{x}_{\boldsymbol{\nu}(t)}, \boldsymbol{\nu}(t))) \\
&= \int_0^1 \left\langle \mathbf{f}_{\theta}(\mathbf{x}_0, \boldsymbol{\nu}(t)), \frac{d}{dt} \boldsymbol{\nu}(t) \right\rangle dt \tag{30}
\end{aligned}$$

313 Thus $\mathcal{L}_{\text{diffusion}}$ can be interpreted as the amount of work done along the trajectory $\boldsymbol{\nu}(0) \rightarrow \boldsymbol{\nu}(1)$ in the
314 presence of a vector field $\mathbf{f}_{\theta}(\mathbf{x}_0, \boldsymbol{\nu}(\mathbf{z}, t))$. From the perspective of thermodynamics, this is precisely
315 equal to the amount of heat lost into the environment during the process of transition between 2
316 equilibria via the noise schedule specified by $\boldsymbol{\nu}(t)$.

317 B.5 Log likelihood and Noise Schedules: A Thermodynamics perspective

318 A diffusion model characterizes a quasi-static process that occurs between two equilibrium distri-
319 butions: $q(\mathbf{x}_0) \rightarrow q(\mathbf{x}_1)$, via a stochastic trajectory (Sohl-Dickstein et al., 2015). According to
320 Spinney & Ford (2012), it is demonstrated that the diffusion schedule or the noising process plays
321 a pivotal role in determining the "measure of irreversibility" for this stochastic trajectory which is
322 expressed as $\log \frac{P_F(\mathbf{x}_{0:1})}{P_B(\mathbf{x}_{1:0})}$. $P_F(\mathbf{x}_{0:1})$ represents the probability of observing the forward path $\mathbf{x}_{0:1}$
323 and $P_B(\mathbf{x}_{1:0})$ represents the probability of observing the reverse path $\mathbf{x}_{1:0}$. It's worth noting that

324 $\log \frac{P_F(\mathbf{x}_{0:1})}{P_B(\mathbf{x}_{1:0})}$ corresponds precisely to the ELBO Eq. 1 that we optimize when training a diffusion
 325 model. Consequently, thermodynamics asserts that the noise schedule indeed has an impact on the
 326 log-likelihood of the diffusion model which contradicts Kingma et al. (2021).

327 C Multivariate noise schedule conditioned on context

328 Let’s say we have a context variable $\mathbf{c} \in \mathbb{R}^m$ that captures high level information about \mathbf{x}_0 .
 329 $\alpha_t(\mathbf{c}), \sigma_t(\mathbf{c}) \in \mathbb{R}^{n \times n}$ are diagonal matrices. The timesteps s, t satisfy $0 \leq s < t \leq 1$. Furthermore,
 330 we use the following notations:

$$\alpha_{t|s}(\mathbf{c}) = \frac{\alpha_t(\mathbf{c})}{\alpha_s(\mathbf{c})} \quad (31)$$

$$\sigma_{t|s}^2(\mathbf{c}) = \sigma_t^2(\mathbf{c}) - \frac{\alpha_{t|s}^2(\mathbf{c})}{\sigma_s^2(\mathbf{c})} \quad (32)$$

331 The forward process for such a method is given as:

$$q_\phi(\mathbf{x}_t | \mathbf{x}_s, \mathbf{c}) = \mathcal{N}\left(\alpha_{t|s}(\mathbf{c})\mathbf{x}_s, \sigma_{t|s}^2(\mathbf{c})\right) \quad (33)$$

332 The distribution of \mathbf{x}_t given \mathbf{x}_s is given by (the derivation is similar to Hooeboom & Salimans
 333 (2022)):

$$\begin{aligned} & q_\phi(\mathbf{x}_s | \mathbf{x}_t, \mathbf{x}_0, \mathbf{c}) \\ &= \mathcal{N}\left(\boldsymbol{\mu}_q = \frac{\alpha_{t|s}(\mathbf{c})\sigma_s^2(\mathbf{c})}{\sigma_t^2(\mathbf{c})}\mathbf{x}_t + \frac{\sigma_{t|s}^2(\mathbf{c})\alpha_s(\mathbf{c})}{\sigma_t^2(\mathbf{c})}\mathbf{x}_0, \boldsymbol{\Sigma}_q = \text{diag}\left(\frac{\sigma_s^2(\mathbf{c})\sigma_{t|s}^2(\mathbf{c})}{\sigma_t^2(\mathbf{c})}\right)\right) \end{aligned} \quad (34)$$

334 We explore various parameterizations for $\gamma_\phi(\mathbf{c}, t)$. These schedules are designed in a manner that
 335 guarantees $\gamma_\phi(\mathbf{c}, 0) = \gamma_{\min}\mathbf{1}_d$ and $\gamma_\phi(\mathbf{c}, 1) = \gamma_{\max}\mathbf{1}_d$. Below, we list these parameterizations. The
 336 polynomial parameterization is novel to our work and yields significant performance gains.

337 **Monotonic Neural Network.** (Kingma et al., 2021) We use the monotonic neural network $\gamma_{\text{vdm}}(t)$,
 338 proposed in VDM to express γ as a function of t such that $\gamma_{\text{vdm}}(t) : [0, 1] \rightarrow [\gamma_{\min}, \gamma_{\max}]^d$. Then we
 339 use FiLM conditioning (Perez et al., 2018) in the intermediate layers of this network via a neural
 340 network that maps \mathbf{z} . The activations of the FiLM layer are constrained to be positive.

341 **Sigmoid.** (Ours) We express $\gamma_\phi(\mathbf{c}, t)$ as a sigmoid function in t such that:

342 $\gamma_\phi(\mathbf{c}, t) = \gamma_{\min} + (\gamma_{\max} - \gamma_{\min}) \frac{\text{sigmoid}(\mathbf{a}_\phi(\mathbf{c})t + \mathbf{b}_\phi(\mathbf{c})) - \text{sigmoid}(\mathbf{b}_\phi(\mathbf{c}))}{\text{sigmoid}(\mathbf{a}_\phi(\mathbf{c}) + \mathbf{b}_\phi(\mathbf{c})) - \text{sigmoid}(\mathbf{b}_\phi(\mathbf{c}))}$ where the coefficients $\mathbf{a}_\phi, \mathbf{b}_\phi$ are
 343 parameterized by a neural network such that $\mathbf{a}_\phi : \mathbb{R}^m \rightarrow \mathbb{R}_{>0}^d, \mathbf{b}_\phi : \mathbb{R}^m \rightarrow \mathbb{R}^d$.

344 **Polynomial.** (Ours) We express $\gamma_\phi(\mathbf{c}, t)$ as a monotonic degree 5 polynomial in t . Details about the
 345 exact functional form of this polynomial and its implementation can be found in Suppl. D.2.

346 C.1 context is available during the inference time.

347 Even though \mathbf{c} represents the input \mathbf{x}_0 , it could be available during during inference. For example
 348 \mathbf{c} could be class labels (Dhariwal & Nichol, 2021) or preexisting embeddings from an auto-encoder
 349 (Preechakul et al., 2022).

350 C.1.1 Reverse Process: Approximate

351 Let $\mathbf{x}_\theta(\mathbf{x}_t, \mathbf{c}, t)$ be an approximation for \mathbf{x}_0 . Then we get the following reverse process (for brevity
 352 we write $\mathbf{x}_\theta(\mathbf{x}_t, \mathbf{c}, t)$ as \mathbf{x}_θ):

$$p_\theta(\mathbf{x}_s | \mathbf{x}_t, \mathbf{c}) = \mathcal{N}\left(\boldsymbol{\mu}_p = \frac{\alpha_{t|s}(\mathbf{c})\sigma_s^2(\mathbf{c})}{\sigma_t^2(\mathbf{c})}\mathbf{x}_t + \frac{\sigma_{t|s}^2(\mathbf{c})\alpha_s(\mathbf{c})}{\sigma_t^2(\mathbf{c})}\mathbf{x}_\theta, \boldsymbol{\Sigma}_p = \text{diag}\left(\frac{\sigma_s^2(\mathbf{c})\sigma_{t|s}^2(\mathbf{c})}{\sigma_t^2(\mathbf{c})}\right)\right) \quad (35)$$

353 C.1.2 Diffusion Loss

354 Similar to the derivation of multi-variate $\mathcal{L}_{\text{diffusion}}$ in Eq. 27 we can derive $\mathcal{L}_{\text{diffusion}}$ for this case too:

$$\mathcal{L}_{\text{diffusion}} = -\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n), t \sim U[0,1]} \left[(\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, \mathbf{c}, t))^\top \text{diag} \left(\frac{d}{dt} \boldsymbol{\nu}(t) \right) (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, \mathbf{c}, t)) \right] \quad (36)$$

355 C.1.3 Limitations of this method

356 This approach is very limited where the diffusion process is only conditioned on class labels. Using
 357 pre-existing embeddings like Diff-AE (Preechakul et al., 2022) is also not possible in general and is
 358 only limited to tasks such as attribute manipulation in datasets.

359 C.2 context isn't available during the inference time.

360 If the context, \mathbf{c} is an explicit function of the input \mathbf{x}_0 things become challenging because \mathbf{x}_0 isn't
 361 available during the inference stage. For this reason, Eq. 34 can't be used to parameterize $\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p$ in
 362 $p_\theta(\mathbf{x}_s | \mathbf{x}_t)$. Let $p_\theta(\mathbf{x}_s | \mathbf{x}_t) = \mathcal{N}(\boldsymbol{\mu}_p(\mathbf{x}_t, t), \boldsymbol{\Sigma}_p(\mathbf{x}_t, t))$ where $\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p$ are parameterized directly by a
 363 neural network. Using Eq. 2 we get the following diffusion loss:

$$\begin{aligned} \mathcal{L}_{\text{diffusion}} &= T \mathbb{E}_{i \sim U[0, T]} \text{D}_{\text{KL}}(q(\mathbf{x}_{s(i)} | \mathbf{x}_{t(i)}, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{s(i)} | \mathbf{x}_{t(i)})) \\ &= \mathbb{E}_{q_\phi} \left(\underbrace{\frac{T}{2} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^\top \boldsymbol{\Sigma}_\theta^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)}_{\text{term 1}} + \underbrace{\frac{T}{2} \left(\text{tr}(\boldsymbol{\Sigma}_q \boldsymbol{\Sigma}_p^{-1} - \mathbf{I}_n) - \log \frac{|\boldsymbol{\Sigma}_q|}{|\boldsymbol{\Sigma}_p|} \right)}_{\text{term 2}} \right) \end{aligned} \quad (37)$$

364 C.2.1 Reverse Process: Approximate

365 Due to the challenges associated with parameterizing $\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p$ directly using a neural network we
 366 parameterize \mathbf{c} using a neural network that approximates \mathbf{c} in the reverse process. Let $\mathbf{x}_\theta(\mathbf{x}_t, t)$ be an
 367 approximation for \mathbf{x}_0 . Then we get the following reverse Rrocess (for brevity we write $\mathbf{x}_\theta(\mathbf{x}_t, t)$ as
 368 \mathbf{x}_θ , and \mathbf{c}_θ denotes an approximation to \mathbf{c} in the reverse process.):

$$\begin{aligned} p_\theta(\mathbf{x}_s | \mathbf{x}_t) &= \mathcal{N} \left(\boldsymbol{\mu}_p = \frac{\boldsymbol{\alpha}_{t|s}(\mathbf{c}_\theta) \boldsymbol{\sigma}_s^2(\mathbf{c}_\theta)}{\boldsymbol{\sigma}_t^2(\mathbf{c}_\theta)} \mathbf{x}_t + \frac{\boldsymbol{\sigma}_{t|s}^2(\mathbf{c}_\theta) \boldsymbol{\alpha}_s(\mathbf{c}_\theta)}{\boldsymbol{\sigma}_t^2(\mathbf{c}_\theta)} \mathbf{x}_\theta, \boldsymbol{\Sigma}_p = \text{diag} \left(\frac{\boldsymbol{\sigma}_s^2(\mathbf{c}_\theta) \boldsymbol{\sigma}_{t|s}^2(\mathbf{c}_\theta)}{\boldsymbol{\sigma}_t^2(\mathbf{c}_\theta)} \right) \right) \end{aligned} \quad (38)$$

369 Consider the limiting case where $T \rightarrow \infty$. Let's analyze the 2 terms in Eq. 37 separately.

370 Using Eq. 2 and Eq. 4, **term 1** in Eq. 37 simplifies in the following manner:

$$\begin{aligned} &\lim_{T \rightarrow \infty} \frac{T}{2} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^\top \boldsymbol{\Sigma}_\theta^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) \\ &\lim_{T \rightarrow \infty} \frac{T}{2} \sum_{i=1}^d \frac{((\boldsymbol{\mu}_q)_i - (\boldsymbol{\mu}_p)_i)^2}{(\boldsymbol{\Sigma}_\theta)_i} \end{aligned} \quad (39)$$

Substituting $1/T$ as δ

$$\begin{aligned} &\lim_{\delta \rightarrow 0^+} \sum_{i=1}^d \frac{1}{\delta \boldsymbol{\sigma}_i^2(\mathbf{x}_\theta, t - \delta) \left(1 - \frac{\boldsymbol{\nu}_i(\mathbf{x}_\theta, t)}{\boldsymbol{\nu}_i(\mathbf{x}_\theta, t - \delta)} \right)} \times \\ &\quad \left[\frac{\boldsymbol{\alpha}_i(\mathbf{x}, t - \delta)}{\boldsymbol{\alpha}_i(\mathbf{x}, t)} \frac{\boldsymbol{\nu}_i(\mathbf{x}, t)}{\boldsymbol{\nu}_i(\mathbf{x}, t - \delta)} \mathbf{z}_t + \boldsymbol{\alpha}_i(\mathbf{x}, t - \delta) \left(1 - \frac{\boldsymbol{\nu}_i(\mathbf{x}, t)}{\boldsymbol{\nu}_i(\mathbf{x}, t - \delta)} \right) x_i \right. \\ &\quad \left. - \frac{\boldsymbol{\alpha}_i(\mathbf{x}_\theta, t - \delta)}{\boldsymbol{\alpha}_i(\mathbf{x}_\theta, t)} \frac{\boldsymbol{\nu}_i(\mathbf{x}_\theta, t)}{\boldsymbol{\nu}_i(\mathbf{x}_\theta, t - \delta)} \mathbf{z}_t + \boldsymbol{\alpha}_i(\mathbf{x}_\theta, t - \delta) \left(1 - \frac{\boldsymbol{\nu}_i(\mathbf{x}_\theta, t)}{\boldsymbol{\nu}_i(\mathbf{x}_\theta, t - \delta)} \right) (x_\theta)_i \right]^2 \end{aligned} \quad (40)$$

371 Consider the scalar case: substituting $\delta = 1/T$,

$$\begin{aligned} & \lim_{\delta \rightarrow 0} \frac{1}{\delta \sigma^2(\mathbf{x}_\theta, t - \delta) \left(1 - \frac{\nu(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t - \delta)}\right)} \times \\ & \left[\frac{\alpha(\mathbf{x}, t - \delta)}{\alpha(\mathbf{x}, t)} \frac{\nu(\mathbf{x}, t)}{\nu(\mathbf{x}, t - \delta)} \mathbf{z}_t + \alpha(\mathbf{x}, t - \delta) \left(1 - \frac{\nu(\mathbf{x}, t)}{\nu(\mathbf{x}, t - \delta)}\right) \mathbf{x} \right. \\ & \left. - \frac{\alpha(\mathbf{x}_\theta, t - \delta)}{\alpha(\mathbf{x}_\theta, t)} \frac{\nu(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t - \delta)} \mathbf{z}_t + \alpha(\mathbf{x}_\theta, t - \delta) \left(1 - \frac{\nu(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t - \delta)}\right) \mathbf{x}_\theta \right]^2 \end{aligned} \quad (41)$$

372 Notice that this equation is in indeterminate for when we substitute $\delta = 0$. One can apply L'Hospital
373 rule twice or break it down into 3 terms below. For this reason let's write it as

$$\begin{aligned} \text{expression 1: } & \lim_{\delta \rightarrow 0} \frac{1}{\delta} \times \left[\frac{\alpha(\mathbf{x}, t - \delta)}{\alpha(\mathbf{x}, t)} \frac{\nu(\mathbf{x}, t)}{\nu(\mathbf{x}, t - \delta)} \mathbf{z}_t + \alpha(\mathbf{x}, t - \delta) \left(1 - \frac{\nu(\mathbf{x}, t)}{\nu(\mathbf{x}, t - \delta)}\right) \mathbf{x} \right. \\ & \left. - \frac{\alpha(\mathbf{x}_\theta, t - \delta)}{\alpha(\mathbf{x}_\theta, t)} \frac{\nu(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t - \delta)} \mathbf{z}_t + \alpha(\mathbf{x}_\theta, t - \delta) \left(1 - \frac{\nu(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t - \delta)}\right) \mathbf{x}_\theta \right] \end{aligned} \quad (42)$$

$$\begin{aligned} \text{expression 2: } & \lim_{\delta \rightarrow 0} \frac{1}{\left(1 - \frac{\nu(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t - \delta)}\right)} \times \left[\frac{\alpha(\mathbf{x}, t - \delta)}{\alpha(\mathbf{x}, t)} \frac{\nu(\mathbf{x}, t)}{\nu(\mathbf{x}, t - \delta)} \mathbf{z}_t + \alpha(\mathbf{x}, t - \delta) \left(1 - \frac{\nu(\mathbf{x}, t)}{\nu(\mathbf{x}, t - \delta)}\right) \mathbf{x} \right. \\ & \left. - \frac{\alpha(\mathbf{x}_\theta, t - \delta)}{\alpha(\mathbf{x}_\theta, t)} \frac{\nu(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t - \delta)} \mathbf{z}_t + \alpha(\mathbf{x}_\theta, t - \delta) \left(1 - \frac{\nu(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t - \delta)}\right) \mathbf{x}_\theta \right]^2 \end{aligned} \quad (43)$$

374 Applying L'Hospital rule in expression 1 we get,

$$\begin{aligned} \frac{d}{d\delta} \left(\frac{\alpha(\mathbf{x}, t - \delta)}{\alpha(\mathbf{x}, t)} \frac{\nu(\mathbf{x}, t)}{\nu(\mathbf{x}, t - \delta)} \right) \Bigg|_{\delta=0} &= \frac{\nu(\mathbf{x}, t) - \nu(\mathbf{x}, t) \alpha'(\mathbf{x}, t) + \alpha(\mathbf{x}, t) \nu'(\mathbf{x}, t)}{\alpha(\mathbf{x}, t) \nu^2(\mathbf{x}, t)} \\ &= \frac{-\alpha'(\mathbf{x}, t)}{\alpha(\mathbf{x}, t)} + \frac{\nu'(\mathbf{x}, t)}{\nu(\mathbf{x}, t)} \end{aligned} \quad (44)$$

$$\frac{d}{d\delta} \alpha(\mathbf{x}, t - \delta) \left(1 - \frac{\nu(\mathbf{x}, t)}{\nu(\mathbf{x}, t - \delta)}\right) \Bigg|_{\delta=0} = -\alpha(\mathbf{x}, t) \frac{\nu'(\mathbf{x}, t)}{\nu(\mathbf{x}, t)} \quad (45)$$

$$\left[\left(\frac{-\alpha'(\mathbf{x}, t)}{\alpha(\mathbf{x}, t)} + \frac{\nu'(\mathbf{x}, t)}{\nu(\mathbf{x}, t)} + \frac{\alpha'(\mathbf{x}_\theta, t)}{\alpha(\mathbf{x}_\theta, t)} - \frac{\nu'(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t)} \right) \mathbf{z}_t \right. \quad (46)$$

$$\left. - \alpha(\mathbf{x}, t) \frac{\nu'(\mathbf{x}, t)}{\nu(\mathbf{x}, t)} \mathbf{x} + \alpha(\mathbf{x}_\theta, t) \frac{\nu'(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t)} \mathbf{x}_\theta \right]^2 \times \frac{\nu(\mathbf{x}, t)}{\nu'(\mathbf{x}, t)} \quad (47)$$

375 Thus the final result:

$$\begin{aligned}
& \sum_{i=1}^d \left[\left(\frac{-\alpha_i'(\mathbf{x}, t)}{\alpha_i(\mathbf{x}, t)} + \frac{\nu_i'(\mathbf{x}, t)}{\nu_i(\mathbf{x}, t)} + \frac{\alpha_i'(\mathbf{x}_\theta, t)}{\alpha_i(\mathbf{x}_\theta, t)} - \frac{\nu_i'(\mathbf{x}_\theta, t)}{\nu_i(\mathbf{x}_\theta, t)} \right) \mathbf{z}_t \right. \\
& \quad \left. - \alpha_i(\mathbf{x}, t) \frac{\nu_i'(\mathbf{x}, t)}{\nu_i(\mathbf{x}, t)} \mathbf{x} + \alpha_i(\mathbf{x}_\theta, t) \frac{\nu_i'(\mathbf{x}_\theta, t)}{\nu_i(\mathbf{x}_\theta, t)} \mathbf{x}_\theta \right]^2 \times \frac{\nu_i(\mathbf{x}, t)}{\nu_i'(\mathbf{x}, t)} \\
& = \Lambda^\top \text{diag} \left(\frac{\nu(\mathbf{x}, t)}{\nu'(\mathbf{x}, t)} \right) \Lambda \\
& \text{where } \Lambda = \left[\left(\frac{-\alpha'(\mathbf{x}, t)}{\alpha(\mathbf{x}, t)} + \frac{\nu'(\mathbf{x}, t)}{\nu(\mathbf{x}, t)} + \frac{\alpha'(\mathbf{x}_\theta, t)}{\alpha(\mathbf{x}_\theta, t)} - \frac{\nu'(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t)} \right) \mathbf{z}_t - \alpha(\mathbf{x}, t) \frac{\nu'(\mathbf{x}, t)}{\nu(\mathbf{x}, t)} \mathbf{x} + \alpha(\mathbf{x}_\theta, t) \frac{\nu'(\mathbf{x}_\theta, t)}{\nu(\mathbf{x}_\theta, t)} \mathbf{x}_\theta \right] \quad (48)
\end{aligned}$$

376 For the second term we have the following:

$$\begin{aligned}
& \lim_{T \rightarrow \infty} \frac{T}{2} \left(\text{tr}(\Sigma_q \Sigma_p^{-1} - \mathbf{I}_n) - \log \frac{|\Sigma_q|}{|\Sigma_p|} \right) \\
& = \lim_{T \rightarrow \infty} \frac{T}{2} \left[\text{tr} \left(\text{diag} \left(\sigma^2(\mathbf{c}, s) \left(1 - \frac{\nu(\mathbf{c}, t)}{\nu(\mathbf{c}, s)} \right) \right) / \text{diag} \left(\sigma^2(\mathbf{c}_\theta, s) \left(1 - \frac{\nu(\mathbf{c}_\theta, t)}{\nu(\mathbf{c}_\theta, s)} \right) \right) - \mathbf{I}_n \right) \right. \\
& \quad \left. - \log \frac{\left| \text{diag} \left(\sigma^2(\mathbf{c}, s) \left(1 - \frac{\nu(\mathbf{c}, t)}{\nu(\mathbf{c}, s)} \right) \right) \right|}{\left| \text{diag} \left(\sigma^2(\mathbf{c}_\theta, s) \left(1 - \frac{\nu(\mathbf{c}_\theta, t)}{\nu(\mathbf{c}_\theta, s)} \right) \right) \right|} \right] \\
& = \lim_{T \rightarrow \infty} \frac{T}{2} \sum_{i=1}^d \left(\frac{\sigma_i^2(\mathbf{c}, s) \left(1 - \frac{\nu_i(\mathbf{c}, t)}{\nu_i(\mathbf{c}, s)} \right)}{\sigma_i^2(\mathbf{c}_\theta, s) \left(1 - \frac{\nu_i(\mathbf{c}_\theta, t)}{\nu_i(\mathbf{c}_\theta, s)} \right)} - 1 - \log \frac{\sigma_i^2(\mathbf{c}, s) \left(1 - \frac{\nu_i(\mathbf{c}, t)}{\nu_i(\mathbf{c}, s)} \right)}{\sigma_i^2(\mathbf{c}_\theta, s) \left(1 - \frac{\nu_i(\mathbf{c}_\theta, t)}{\nu_i(\mathbf{c}_\theta, s)} \right)} \right) \quad (49)
\end{aligned}$$

(50)

$$377 \quad \text{Let } p_i = \frac{\sigma_i^2(\mathbf{c}, s) \left(1 - \frac{\nu_i(\mathbf{c}, t)}{\nu_i(\mathbf{c}, s)} \right)}{\sigma_i^2(\mathbf{c}_\theta, s) \left(1 - \frac{\nu_i(\mathbf{c}_\theta, t)}{\nu_i(\mathbf{c}_\theta, s)} \right)}$$

378 The sequence $\lim_{T \rightarrow \infty} \frac{T}{2} \sum_{i=1}^d (p_i - 1 - \log p_i)$ converges iff $\lim_{T \rightarrow \infty} \sum_{i=1}^d (p_i - 1 - \log p_i) = 0$.

379 Notice that the function $f(x) = x - 1 - \log x \geq 0 \quad \forall x \in \mathbb{R}$ and the equality holds for $x = 1$. Thus,

380 the condition $\lim_{T \rightarrow \infty} \frac{T}{2} \sum_{i=1}^d (p_i - 1 - \log p_i)$ holds iff $\lim_{T \rightarrow \infty} p_i = 0 \quad \forall i \in \{1, \dots, d\}$. Thus,

$$\begin{aligned} \lim_{T \rightarrow \infty} p_i &= 1 \\ \implies \lim_{T \rightarrow \infty} \left(\frac{\sigma_i^2(\mathbf{c}, s) \left(1 - \frac{\nu_i(\mathbf{c}, t)}{\nu_i(\mathbf{c}, s)}\right)}{\sigma_i^2(\mathbf{c}_\theta, s) \left(1 - \frac{\nu_i(\mathbf{c}_\theta, t)}{\nu_i(\mathbf{c}_\theta, s)}\right)} \right) &= 1 \end{aligned}$$

Substituting $1/T$ as δ ,

$$\begin{aligned} \implies \lim_{\delta \rightarrow 0^+} \left(\frac{\sigma_i^2(\mathbf{c}, t - \delta) \left(1 - \frac{\nu_i(\mathbf{c}, t)}{\nu_i(\mathbf{c}, t - \delta)}\right)}{\sigma_i^2(\mathbf{c}_\theta, t - \delta) \left(1 - \frac{\nu_i(\mathbf{c}_\theta, t)}{\nu_i(\mathbf{c}_\theta, t - \delta)}\right)} \right) &= 1 \\ \implies \frac{\sigma_i^2(\mathbf{c}, t)}{\sigma_i^2(\mathbf{c}_\theta, t)} \lim_{\delta \rightarrow 0^+} \left(\frac{1 - \frac{\nu_i(\mathbf{c}, t)}{\nu_i(\mathbf{c}, t - \delta)}}{1 - \frac{\nu_i(\mathbf{c}_\theta, t)}{\nu_i(\mathbf{c}_\theta, t - \delta)}} \right) &= 1 \end{aligned}$$

Applying L'Hospital rule,

$$\begin{aligned} \implies \frac{\sigma_i^2(\mathbf{c}, t)}{\sigma_i^2(\mathbf{c}_\theta, t)} \left(\frac{-\nu_i'(\mathbf{c}, t)}{\nu_i(\mathbf{c}, t)} \right) &= 1 \\ \implies \frac{\sigma_i^2(\mathbf{c}, t)}{\sigma_i^2(\mathbf{c}_\theta, t)} \left(\frac{\nu_i'(\mathbf{c}, t)\nu_i(\mathbf{c}_\theta, t)}{\nu_i(\mathbf{c}, t)\nu_i'(\mathbf{c}_\theta, t)} \right) &= 1 \end{aligned} \tag{51}$$

381 In the vector form the above equation can be written as,

$$\frac{\sigma_t^2(\mathbf{c})\boldsymbol{\nu}_t(\mathbf{c}_\theta)\nabla_t\boldsymbol{\nu}(\mathbf{c}, t)}{\sigma_t^2(\mathbf{c}_\theta)\boldsymbol{\nu}_t(\mathbf{c})\nabla_t\boldsymbol{\nu}(\mathbf{c}_\theta, t)} \rightarrow \mathbf{1}_d \tag{52}$$

382 Eq. 52 holds if:

- 383 • $x_\theta = x_0$ i.e. the unet can perfectly map \mathbf{x}_t to $\mathbf{x}_0 \forall t \in [0, 1]$ which is unrealistic.
- 384 • Clever parameterizations for σ, α, ν that ensure Eq. 52 holds.

385 Because of aforementioned challenges we evaluate this method with finite $T = 1000$. We demonstrate
386 the performance of the model empirically in Fig. 1.

387 C.2.2 Recovering VDM

388 If we substitute $\boldsymbol{\nu}_t(\mathbf{c}), \boldsymbol{\nu}_t(\mathbf{c}_\theta)$ with $\boldsymbol{\nu}(t)$ (since the SNR isn't conditioned on the context \mathbf{c}),
389 $\sigma_t(\mathbf{c}_\theta), \sigma_t(\mathbf{c})$ with σ_t and $\alpha_t(\mathbf{c}_\theta), \alpha_t(\mathbf{c})$ with α_t , Eq. 39 reduces to the intermediate loss in VDM
390 i.e. $\frac{1}{2}(\mathbf{x}_\theta - \mathbf{x}_0)^\top (\nabla_t \boldsymbol{\nu}(t)) (\mathbf{x}_\theta - \mathbf{x}_0)$ and Eq. 49 reduces to 0.

391 C.3 Challenges in Conditioning on Context

392 Note that the model $p_\theta(\mathbf{x}_{0:1}|\mathbf{c})$ implicitly assumes the availability of \mathbf{c} at generation time. Sometimes,
393 this context may be available, such as when we condition on a label. We may then fit a conditional
394 diffusion process with a standard diffusion objective $\mathbb{E}_{\mathbf{x}_0, \mathbf{c}}[\text{ELBO}(\mathbf{x}_0, p_\theta(\mathbf{x}_{0:1}|\mathbf{c}), q_\phi(\mathbf{x}_{0:1}|\mathbf{c}))]$, in
395 which both the forward and the backward processes are conditioned on \mathbf{c} (see Sec. 2.3).

396 When \mathbf{c} is not known at generation time, we may fit a model p_θ that does not condition on \mathbf{c} . Unfortu-
397 nately, this also forces us to define $p_\theta(\mathbf{x}_s|\mathbf{x}_t) = \mathcal{N}(\boldsymbol{\mu}_p(\mathbf{x}_t, t), \boldsymbol{\Sigma}_p(\mathbf{x}_t, t))$ where $\boldsymbol{\mu}_p(\mathbf{x}_t, t), \boldsymbol{\Sigma}_p(\mathbf{x}_t, t)$
398 is parameterized directly by a neural network. We can no longer use a noise parameterization
399 $\epsilon_\theta(\mathbf{x}_t, t) = (\mathbf{x}_t - \alpha_t(\mathbf{c})\mathbf{x}_\theta(\mathbf{x}_t, t, \mathbf{c}))/\sigma_t(\mathbf{c})$ because it requires us to compute $\alpha_t(\mathbf{c})$ and $\sigma_t(\mathbf{c})$,
400 which we do not know. Since noise parameterization plays a key role in the sample quality of
401 diffusion models (Ho et al., 2020), this approach limits performance.

402 The other approach is to approximate \mathbf{c} using a neural network, $\mathbf{c}_\theta(\mathbf{x}_t, t)$. This would allow us
403 to write $p_\theta(\mathbf{x}_s|\mathbf{x}_t) = q_\phi(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0 = \mathbf{x}_\theta(\mathbf{x}_t, t), \mathbf{c} = \mathbf{c}_\theta(\mathbf{x}_t, t))$. Unfortunately, this introduces
404 instability in the learning objective, which we observe both theoretically and empirically. Specifically,

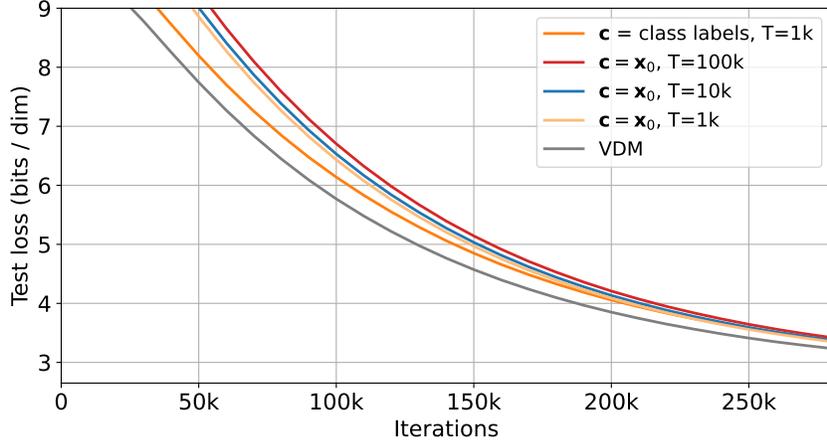


Figure 1: For \mathbf{c} = “class labels” or $\mathbf{c} = \mathbf{x}_0$ the likelihood estimates are worse than VDM. For $\mathbf{c} = \mathbf{x}_0$ we see that the VLB degrades with increasing T whereas for VDM and MULAN it improves for increasing T ; see Kingma et al. (2021). This empirical observation is consistent with our mathematical insights earlier. As these models consistently exhibit inferior performance w.r.t VDM, in line with our initial conjectures, we refrain from training them beyond 300k iterations due to the substantial computational cost involved.

405 in Suppl. C we show that the learning objective diverges unless the following condition holds true:
 406 $\lim_{T \rightarrow \infty} T \frac{\sigma_t^2(\mathbf{x}_0) \nu_t(\mathbf{x}_0) \nu_t'(\mathbf{x}_\theta)}{\sigma_t^2(\mathbf{x}_\theta) \nu_t(\mathbf{x}_\theta) \nu_t'(\mathbf{x}_0)} \rightarrow \mathbf{I}_d$ pointwise across t . Experiments in Suppl. C.4 confirm this issue.

407 C.4 Experimental results

408 In Fig. 1 we demonstrate that the multivariate diffusion processes where \mathbf{c} = “class labels” or $\mathbf{c} = \mathbf{x}_0$
 409 perform worse than VDM. Since a continuous time formulation i.e. $T \rightarrow \infty$ for the case when $\mathbf{c} = \mathbf{x}_0$
 410 isn’t possible (unlike MULAN or VDM) we evaluate these models in the discrete time setting where
 411 we use $T = 1000$. Furthermore we also ablate $T = 10k, 100k$ for $\mathbf{c} = \mathbf{x}_0$ to show that the VLB
 412 degrades with increasing T whereas for VDM and MULAN it improves for increasing T ; see Kingma
 413 et al. (2021). This empirical observation is consistent with our mathematical insights earlier. As these
 414 models consistently exhibit inferior performance w.r.t VDM, in line with our initial conjectures, we
 415 refrain from training them beyond 300k iterations due to the substantial computational cost involved.

416 D MULAN: MULTIVARIATE Latent Auxiliary variable Noise Schedule

417 D.1 Parameterization in the reverse process

418 D.1.1 Noise parameterization

419 Since the forward pass is given by $\mathbf{x}_t = \alpha_t(\mathbf{z})\mathbf{x}_0 + \sigma_t(\mathbf{z})\epsilon_t$, we can write the noise ϵ_t in terms of
 420 $\mathbf{x}_0, \mathbf{x}_t$ in the following manner:

$$\epsilon_t = \frac{\mathbf{x}_t - \alpha_t(\mathbf{z})\mathbf{x}_0}{\sigma_t(\mathbf{z})} \quad (53)$$

421 Following Dhariwal & Nichol (2021); Kingma et al. (2021), instead of parameterizing $\mathbf{x}_\theta(\mathbf{x}_t, \mathbf{z}, t)$
 422 using a neural network, we use Eq. 53 to parameterize the denoising model in terms of a noise
 423 prediction model $\epsilon_\theta(\mathbf{x}_t, \mathbf{z}, t)$,

$$\epsilon_\theta(\mathbf{x}_t, \mathbf{z}, t) = \frac{\mathbf{x}_t - \alpha_t(\mathbf{z})\mathbf{x}_\theta(\mathbf{x}_t, \mathbf{z}, t)}{\sigma_t(\mathbf{z})} \quad (54)$$

Table 2: Likelihood in bits per dimension (BPD) (mean and 95% confidence interval), on the test set of CIFAR-10 computed using VLB estimate.

parameterization	Num training steps	CIFAR-10 (\downarrow)
Noise parameterization	10M	2.60 ± 10^{-3}
Velocity parameterization	8M	2.59 ± 10^{-3}

424 D.1.2 Velocity parameterization

425 Following [Salimans & Ho \(2022\)](#); [Zheng et al. \(2023\)](#), we explore another parameterization of the
426 denoising network which is given by

$$\mathbf{v}_\theta(\mathbf{x}_t, \mathbf{z}, t) = \frac{\boldsymbol{\alpha}_t(\mathbf{z})\mathbf{x}_t - \mathbf{x}_\theta(\mathbf{x}_t, \mathbf{z}, t)}{\boldsymbol{\sigma}_t(\mathbf{z})} \quad (55)$$

427 In practice, Velocity parameterization leads to a better performance than noise parameterization; as
428 illustrated in [Table 2](#).

429 D.2 Polynomial Noise Schedule

Let $f(x; \psi)$ be a scalar-valued polynomial of degree n with coefficients $\psi \in \mathbb{R}^{n+1}$ expressed as:

$$f(x; \psi) = \psi_n x^n + \psi_{n-1} x^{n-1} + \dots + \psi_1 x + \psi_0,$$

430 and denote its derivative with respect to x as $\frac{d}{dx}f(x; \psi)$, represented by $f'(x; \psi)$. Now we'd like to
431 find least n such that $f(x; \psi)$ satisfies the following properties:

- 432 1. $f(x; \psi)$ is monotonically increasing, i.e. $f'(x; \psi) \geq 0 \forall x \in \mathbb{R}, \psi \in \mathbb{R}^{n+1}$.
- 433 2. $f'(x_1; \psi) = 0, f'(x_2; \psi) = 0 \exists x_1, x_2 \in \mathbb{C}, x_1 \neq x_2, \forall \psi \in \mathbb{R}^{n+1}$.

434 For the **first** condition to hold, we can design $f'(x; \psi)$ such that it's a perfect square with real /
435 imaginary roots. That way $f'(x; \psi) \geq 0 \forall x \in \mathbb{R}, \psi \in \mathbb{R}^{n+1}$. This means that $f'(x; \psi)$ is an even
436 degree polynomial, i.e. the degree of $f'(x; \psi)$ can take the following values: 2, 4, ... Also, note that
437 at least half of the roots of $f'(x; \psi)$ are repeated since $f'(x; \psi)$ can be expressed as a perfect square,
438 i.e., if $f'(x; \psi)$ has a degree 2 then it has exactly 1 unique root (real / imaginary), if $f'(x; \psi)$ has a
439 degree 4 then it has at most 2 unique roots (real / imaginary), and so on.

440 For the **second** condition to hold, $f'(x; \psi)$ needs to have at least 2 unique roots $\exists \psi \in \mathbb{R}^{n+1}$. For this
441 reason $f'(x; \psi)$ is a polynomial of degree 4. Thus, $f'(x; \psi)$ can be written as $f'(x; \psi) = (\psi_3 x^2 +$
442 $\psi_2 x + \psi_1)^2$. This ensures that $\exists \psi \in \mathbb{R}^5$ s.t. $f'(x; \psi) = 0$ twice in $x \in \mathbb{R}$, and $f'(x; \psi) \geq 0 \forall \psi \in \mathbb{R}^5$.

443 Thus, $f(x; \psi)$ takes the following functional form:

$$\begin{aligned} f(x; \psi) &= \int (\psi_3 x^2 + \psi_2 x + \psi_1)^2 dx \\ &= \frac{\psi_3^2}{5} x^5 + \frac{\psi_3 \psi_2}{2} x^4 + \frac{\psi_2^2 + 2\psi_3 \psi_1}{3} x^3 + \psi_2 \psi_1 x^2 + \psi_1^2 x + \text{constant}. \end{aligned} \quad (56)$$

For the above-mentioned reasons we express $\gamma(\mathbf{c}, t) : \mathbb{R}^m \times [0, 1] \rightarrow \mathbb{R}^d$ as a degree 5 polynomial in
 t . We define neural networks $\mathbf{a}_\phi(\mathbf{c}) : \mathbb{R}^m \rightarrow \mathbb{R}^d$, $\mathbf{b}_\phi(\mathbf{c}) : \mathbb{R}^m \rightarrow \mathbb{R}^d$, and $\mathbf{d}_\phi(\mathbf{c}) : \mathbb{R}^m \rightarrow \mathbb{R}^d$ with
parameters ϕ . Let $f_\phi : \mathbb{R}^m \times [0, 1] \rightarrow \mathbb{R}^d$ be defined as:

$$f_\phi(\mathbf{c}, t) = \frac{\mathbf{a}_\phi^2(\mathbf{c})}{5} t^5 + \frac{\mathbf{a}_\phi(\mathbf{c})\mathbf{b}_\phi(\mathbf{c})}{2} t^4 + \frac{\mathbf{b}_\phi^2(\mathbf{c}) + 2\mathbf{a}_\phi(\mathbf{c})\mathbf{d}_\phi(\mathbf{c})}{3} t^3 + \mathbf{b}_\phi(\mathbf{c})\mathbf{d}_\phi(\mathbf{c})t^2 + \mathbf{d}_\phi^2(\mathbf{c})t$$

444 where the multiplication and division operations are elementwise. The the noise schedule, $\gamma(\mathbf{c}, t)$, is
445 given as follows:

$$\gamma_\phi(\mathbf{c}, t) = \gamma_{\min} + (\gamma_{\max} - \gamma_{\min}) \frac{f_\phi(\mathbf{c}, t)}{f_\phi(\mathbf{c}, t=1)} \quad (57)$$

446 Notice that $\gamma_\phi(\mathbf{c}, t)$ has these interesting properties:

- 447 • Is increasing in $t \in [0, 1]$ which is crucial as mentioned in Sec. D.4.
- 448 • $\gamma_\phi(\mathbf{c}, t)$ has end points at $t = 0$ and $t = 1$ which the user can specify via γ_{\min} and γ_{\max} .
449 Specifically, $\gamma_\phi(\mathbf{c}, t = 0) = \gamma_{\min}$ and $\gamma_\phi(\mathbf{c}, t = 1) = \gamma_{\max}$.
- 450 • Its time-derivative i.e. $\nabla_t \gamma_\phi(\mathbf{c}, t)$ **can** be zero twice in $t \in [0, 1]$. This isn't a necessary
451 condition but it's nice to have a flexible noise schedule whose time-derivative can be 0 at the
452 beginning and the end of the diffusion process.

453 D.3 Variational Lower Bound

454 In this section we derive the VLB. For ease of reading we use the notation \mathbf{x}_t to denote $\mathbf{x}_{t(i)}$ and
455 \mathbf{x}_{t-1} to denote $\mathbf{x}_{t(i-1)} \equiv \mathbf{x}_{s(i)}$ in the following derivation.

$$\begin{aligned}
& -\log p_\theta(\mathbf{x}_0) \\
& \leq \mathbb{E}_{q_\phi} \left[-\log \frac{p_\theta(\mathbf{z}, \mathbf{x}_{0:T})}{q_\phi(\mathbf{z}, \mathbf{x}_{1:T} | \mathbf{x}_0)} \right] \\
& = \mathbb{E}_{q_\phi} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T-1} | \mathbf{z}, \mathbf{x}_T)}{q_\phi(\mathbf{z}, \mathbf{x}_{1:T} | \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_T) - \log p_\theta(\mathbf{z}) \right] \\
& = \mathbb{E}_{q_\phi} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T-1} | \mathbf{z}, \mathbf{x}_T)}{q_\phi(\mathbf{x}_{1:T} | \mathbf{z}, \mathbf{x}_0)} - \log \frac{1}{q_\phi(\mathbf{z} | \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_T) - \log p_\theta(\mathbf{z}) \right] \\
& = \mathbb{E}_{q_\phi} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T-1} | \mathbf{z}, \mathbf{x}_T)}{q_\phi(\mathbf{x}_{1:T} | \mathbf{z}, \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_T) - \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x}_0)} \right] \\
& = \mathbb{E}_{q_\phi} \left[-\sum_{t=1}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_t)}{q_\phi(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{z}, \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_T) - \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x}_0)} \right] \\
& = \mathbb{E}_{q_\phi} \left[-\log \frac{p_\theta(\mathbf{x}_0 | \mathbf{z}, \mathbf{x}_1)}{q_\phi(\mathbf{x}_1 | \mathbf{x}_0, \mathbf{z})} - \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_t)}{q_\phi(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{z}, \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_T) - \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x}_0)} \right] \\
& = \mathbb{E}_{q_\phi} \left[-\log \frac{p_\theta(\mathbf{x}_0 | \mathbf{z}, \mathbf{x}_1)}{q_\phi(\mathbf{x}_1 | \mathbf{x}_0, \mathbf{z})} - \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_t) q_\phi(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_0)}{q_\phi(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{z}, \mathbf{x}_0) q_\phi(\mathbf{x}_t | \mathbf{z}, \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_T) - \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x}_0)} \right] \\
& = \mathbb{E}_{q_\phi} \left[-\log \frac{p_\theta(\mathbf{x}_0 | \mathbf{z}, \mathbf{x}_1)}{q_\phi(\mathbf{x}_1 | \mathbf{x}_0, \mathbf{z})} - \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_t)}{q_\phi(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{z}, \mathbf{x}_0)} - \sum_{t=2}^T \log \frac{q_\phi(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_0)}{q_\phi(\mathbf{x}_t | \mathbf{z}, \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_T) - \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x}_0)} \right] \\
& = \mathbb{E}_{q_\phi} \left[-\log \frac{p_\theta(\mathbf{x}_0 | \mathbf{z}, \mathbf{x}_1)}{q_\phi(\mathbf{x}_1 | \mathbf{x}_0, \mathbf{z})} - \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_t)}{q_\phi(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{z}, \mathbf{x}_0)} - \log \frac{q(\mathbf{x}_1 | \mathbf{z}, \mathbf{x}_0)}{q_\phi(\mathbf{x}_T | \mathbf{z}, \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_T) - \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x}_0)} \right] \\
& = \mathbb{E}_{q_\phi} \left[-\log p_\theta(\mathbf{x}_0 | \mathbf{z}, \mathbf{x}_1) - \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_t)}{q_\phi(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{z}, \mathbf{x}_0)} - \log \frac{1}{q_\phi(\mathbf{x}_T | \mathbf{z}, \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_T) - \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x}_0)} \right] \\
& = \mathbb{E}_{q_\phi} \left[-\log p_\theta(\mathbf{x}_0 | \mathbf{z}, \mathbf{x}_1) - \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_t)}{q_\phi(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{z}, \mathbf{x}_0)} - \log \frac{p_\theta(\mathbf{x}_T)}{q_\phi(\mathbf{x}_T | \mathbf{z}, \mathbf{x}_0)} - \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x}_0)} \right] \\
& = \mathbb{E}_{q_\phi} \left[\underbrace{-\log p_\theta(\mathbf{x}_0 | \mathbf{z}, \mathbf{x}_1)}_{\mathcal{L}_{\text{recons}}} + \underbrace{\sum_{t=2}^T \text{D}_{\text{KL}}[p_\theta(\mathbf{x}_{t-1} | \mathbf{z}, \mathbf{x}_t) \| q_\phi(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{z}, \mathbf{x}_0)]}_{\mathcal{L}_{\text{diffusion}}} \right] \\
& + \mathbb{E}_{q_\phi} \left[\underbrace{\text{D}_{\text{KL}}[p_\theta(\mathbf{x}_T) \| q_\phi(\mathbf{x}_T | \mathbf{z}, \mathbf{x}_0)]}_{\mathcal{L}_{\text{prior}}} + \underbrace{\text{D}_{\text{KL}}[p_\theta(\mathbf{z}) \| q(\mathbf{z} | \mathbf{x}_0)]}_{\mathcal{L}_{\text{latent}}} \right] \tag{58}
\end{aligned}$$

456 Switching back to the notation used throughout the paper, the VLB is given as:

$$\begin{aligned}
& -\log p_\theta(\mathbf{x}_0) \\
& \leq \mathbb{E}_{q_\phi} \left[\underbrace{-\log p_\theta(\mathbf{x}_0|\mathbf{z}, \mathbf{x}_1)}_{\mathcal{L}_{\text{recons}}} + \underbrace{\sum_{i=2}^T \text{D}_{\text{KL}}[p_\theta(\mathbf{x}_{s(i)}|\mathbf{z}, \mathbf{x}_{t(i)})||q_\phi(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}, \mathbf{z}, \mathbf{x}_0)]}_{\mathcal{L}_{\text{diffusion}}} \right] \\
& \quad + \mathbb{E}_{q_\phi} \left[\underbrace{\text{D}_{\text{KL}}[p_\theta(\mathbf{x}_1)||q_\phi(\mathbf{x}_1|\mathbf{z}, \mathbf{x}_0)]}_{\mathcal{L}_{\text{prior}}} + \underbrace{\text{D}_{\text{KL}}[p_\theta(\mathbf{z})||q_\phi(\mathbf{z}|\mathbf{x}_0)]}_{\mathcal{L}_{\text{latent}}} \right] \tag{59}
\end{aligned}$$

457 Next, we derive a precise formula for the learning objective (6) of the auxiliary-variable diffusion
458 model. Using the objective of a diffusion model in (1) we can write (6) as the sum of four terms:

$$\log p_\theta(\mathbf{x}_0) \geq \mathbb{E}_{q_\phi} [\mathcal{L}_{\text{recons}} + \mathcal{L}_{\text{diffusion}} + \mathcal{L}_{\text{prior}} + \mathcal{L}_{\text{latent}}], \tag{60}$$

459 The reconstruction loss, $\mathcal{L}_{\text{recons}}$, can be (stochastically and differentially) estimated using standard
460 techniques; see (Kingma & Welling, 2013), $\mathcal{L}_{\text{prior}} = -\text{D}_{\text{KL}}[q_\phi(\mathbf{x}_1|\mathbf{x}_0, \mathbf{z})||p_\theta(\mathbf{x}_1)]$ is the diffusion
461 prior term, $\mathcal{L}_{\text{latent}} = -\text{D}_{\text{KL}}[q_\phi(\mathbf{z}|\mathbf{x}_0)||p_\theta(\mathbf{z})]$ is the latent prior term, and $\mathcal{L}_{\text{diffusion}}$ is the diffusion loss
462 term, which we examine below. The complete derivation is given in Suppl. D.3.

463 D.3.1 Diffusion Loss

464 **Discrete-Time Diffusion.** We start by defining p_θ in discrete time, and as in Sec. 1, we let $T > 0$
465 be the number of total time steps and define $t(i) = i/T$ and $s(i) = (i-1)/T$ as indexing variables
466 over the time steps. We also use $\mathbf{x}_{0:1}$ to denote the subset of variables associated with these timesteps.
467 Starting with the expression in Eq. 1 and following the steps in Suppl. D, we can write $\mathcal{L}_{\text{diffusion}}$ as:

$$\begin{aligned}
\mathcal{L}_{\text{diffusion}} &= -\sum_{i=2}^T \text{D}_{\text{KL}}[q_\phi(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}, \mathbf{x}_0, \mathbf{z})||p_\theta(\mathbf{x}_{s(i)}|\mathbf{x}_{t(i)}, \mathbf{z})] \\
&= \frac{1}{2} \sum_{i=2}^T [(\epsilon_t - \epsilon_\theta(\mathbf{x}_t, \mathbf{z}, t(i)))^\top \text{diag}(\gamma(\mathbf{z}, s(i)) - \gamma(\mathbf{z}, t(i))) (\epsilon_t - \epsilon_\theta(\mathbf{x}_t, \mathbf{z}, t(i)))] \tag{61}
\end{aligned}$$

468 **Continuous-Time Diffusion.** We can also consider the limit of the above objective as we take an
469 infinitesimally small partition of $t \in [0, 1]$, which corresponds to the limit when $T \rightarrow \infty$. In Suppl. D
470 we show that taking this limit of Eq. 61 yields the continuous-time diffusion loss:

$$\mathcal{L}_{\text{diffusion}} = -\frac{1}{2} \mathbb{E}_{t \sim [0,1]} [(\epsilon_t - \epsilon_\theta(\mathbf{x}_t, \mathbf{z}, t))^\top \text{diag}(\nabla_t \gamma(\mathbf{z}, t)) (\epsilon_t - \epsilon_\theta(\mathbf{x}_t, \mathbf{z}, t))] \tag{62}$$

471 where $\nabla_t \gamma(\mathbf{z}, t) \in \mathbb{R}^d$ denotes the Jacobian of $\gamma(\mathbf{z}, t)$ with respect to the scalar t . We observe that
472 the limit of $T \rightarrow \infty$ yields improved performance, matching the existing theoretical argument by
473 Kingma et al. (2021).

474 D.3.2 Auxiliary latent loss

475 We try two different kinds of priors for $p_\theta(\mathbf{z})$: discrete ($\mathbf{z} \in \{0, 1\}^m$) and continuous ($\mathbf{z} \in \mathbb{R}^m$).

476 **Continuous Auxiliary Latents.** In the case where \mathbf{z} is continuous, we select $p_\theta(\mathbf{z})$ as $\mathcal{N}(\mathbf{0}, \mathbf{I}_m)$.
477 This leads to the following KL loss term:

$$478 \text{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}_0)||p_\theta(\mathbf{z})) = \frac{1}{2}(\boldsymbol{\mu}^\top(\mathbf{x}_0)\boldsymbol{\mu}(\mathbf{x}_0)) + \text{tr}(\boldsymbol{\Sigma}^2(\mathbf{x}_0) - \mathbf{I}_m) - \log |\boldsymbol{\Sigma}^2(\mathbf{x}_0)|.$$

479 **Discrete Auxiliary Latents.** In the case where \mathbf{z} is discrete, we select $p_\theta(\mathbf{z})$ as a uniform dis-
480 tribution. Let $\mathbf{z} \in \{0, 1\}^m$ be a k -hot vector sampled from a discrete Exponential Family dis-
481 tribution $p_\theta(\mathbf{z}; \theta)$ with logits θ . Niepert et al. (2021) show that $\mathbf{z} \sim p_\theta(\mathbf{z}; \theta)$ is equivalent to
482 $\mathbf{z} = \arg \max_{y \in Y} \langle \theta + \epsilon_g, y \rangle$ where ϵ_g denotes the sum of gamma distribution Suppl. E, Y denotes
483 the set of all k -hot vectors of some fixed length m . For $k > 1$, To differentiate through the arg max
484 we use a relaxed estimator, Identity, as proposed by Sahoo et al. (2023). This leads to the following
485 KL loss term: $\text{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}_0)||p_\theta(\mathbf{z})) = -\sum_{i=1}^m q_\phi(\mathbf{z}|\mathbf{x}_0)_i (\log q_\phi(\mathbf{z}|\mathbf{x}_0)_i + \log m)$.

486 **D.4 The Variational Lower Bound as a Line Integral Over The Noise Schedule**

487 Having defined our loss, we now return to the question of whether it is invariant to the choice of
 488 diffusion process. Notice that we may rewrite Eq. 62 in the following vectorized form:

$$\mathcal{L}_{\text{diffusion}} = -\frac{1}{2} \int_0^1 (\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_t, \mathbf{z}, t))^2 \cdot \nabla_t \boldsymbol{\nu}(\mathbf{z}, t) dt \quad (63)$$

489 where the square is applied elementwise. We seek to rewrite (63) as a line integral $\int_a^b \mathbf{f}(\mathbf{r}(t)) \cdot \frac{d}{dt} \mathbf{r}(t) dt$
 490 for some vector field \mathbf{f} and trajectory $\mathbf{r}(t)$. Recall that $\boldsymbol{\nu}(\mathbf{z}, t)$ is monotonically decreasing in each
 491 coordinate as a function of t ; hence, it is invertible on its image, and we can write $t = \boldsymbol{\nu}_z^{-1}(\boldsymbol{\nu}(\mathbf{z}, t))$
 492 for some $\boldsymbol{\nu}_z^{-1}$. Let $\bar{\mathbf{x}}_\theta(\mathbf{x}_{\boldsymbol{\nu}(\mathbf{z}, t)}, \mathbf{z}, \boldsymbol{\nu}(\mathbf{z}, t)) = \mathbf{x}_\theta(\mathbf{x}_{\boldsymbol{\nu}_z^{-1}(\boldsymbol{\nu}(\mathbf{z}, t))}, \mathbf{z}, \boldsymbol{\nu}_z^{-1}(\boldsymbol{\nu}(\mathbf{z}, t)))$ and note that for all t ,
 493 we can write \mathbf{x}_t as $\mathbf{x}_{\boldsymbol{\nu}(\mathbf{z}, t)}$; see Eq. 24, and have $\bar{\mathbf{x}}_\theta(\mathbf{x}_{\boldsymbol{\nu}(\mathbf{z}, t)}, \mathbf{z}, \boldsymbol{\nu}(\mathbf{z}, t)) = \mathbf{x}_\theta(\mathbf{x}_t, \mathbf{z}, t)$. We can then
 494 write the integral in (63) as $\int_0^1 (\mathbf{x}_0 - \bar{\mathbf{x}}_\theta(\mathbf{x}_{\boldsymbol{\nu}(\mathbf{z}, t)}, \mathbf{z}, \boldsymbol{\nu}(\mathbf{z}, t)))^2 \cdot \frac{d}{dt} \boldsymbol{\nu}(\mathbf{z}, t) dt$, which is a line integral
 495 with $\mathbf{f}(\mathbf{r}(t)) \equiv (\mathbf{x}_0 - \bar{\mathbf{x}}_\theta(\mathbf{x}_{\boldsymbol{\nu}(\mathbf{z}, t)}, \mathbf{z}, \boldsymbol{\nu}(\mathbf{z}, t)))^2$ and $\mathbf{r}(t) \equiv \boldsymbol{\nu}(\mathbf{z}, t)$ and .

496 Thus the diffusion loss, $\mathcal{L}_{\text{diffusion}}$, can be interpreted as a measure of work done along the trajectory
 497 $\boldsymbol{\nu}(\mathbf{z}, t)$ in the presence of a vector field \mathbf{f} . Different "trajectories" yield different results for most
 498 integrands, unless its integral corresponds to a conservative force field, which is rarely the case for a
 499 diffusion process (Spinney & Ford, 2012). We empirically observe this in our experiments where
 500 swapping out different multivariate $\boldsymbol{\nu}$ yields different values of the ELBO. In Sec. D.6, we show that
 501 variational diffusion models can be viewed as following only linear trajectories $\boldsymbol{\nu}(t)$, hence their
 502 objective is invariant to the noise schedule. Our method learns a multivariate $\boldsymbol{\nu}$ that yields paths
 503 corresponding to a better ELBO.

504 **D.5 Diffusion Loss**

505 To derive the diffusion loss, $\mathcal{L}_{\text{diffusion}}$ in Eq. 60, we first derive an expression for
 506 $D_{\text{KL}}(q_\phi(\mathbf{x}_s | \mathbf{z}, \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_s | \mathbf{z}, \mathbf{x}_t))$ using Eq. 2 and Eq. 4 in the following manner (details in
 507 Suppl. D):

$$\begin{aligned} & D_{\text{KL}}(q_\phi(\mathbf{x}_s | \mathbf{z}, \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_s | \mathbf{z}, \mathbf{x}_t)) \\ &= \frac{1}{2} \left((\boldsymbol{\mu}_{q_\phi} - \boldsymbol{\mu}_p)^\top \boldsymbol{\Sigma}_\theta^{-1} (\boldsymbol{\mu}_{q_\phi} - \boldsymbol{\mu}_p) + \text{tr}(\boldsymbol{\Sigma}_{q_\phi} \boldsymbol{\Sigma}_p^{-1} - \mathbf{I}_n) - \log \frac{|\boldsymbol{\Sigma}_{q_\phi}|}{|\boldsymbol{\Sigma}_p|} \right) \\ &= \frac{1}{2} \left((\mathbf{x}_0 - \mathbf{x}_\theta)^\top \text{diag}(\boldsymbol{\nu}(\mathbf{z}, s) - \boldsymbol{\nu}(\mathbf{z}, t)) (\mathbf{x}_0 - \mathbf{x}_\theta) \right) \end{aligned} \quad (64)$$

508 Let $\lim_{T \rightarrow \infty} T(\boldsymbol{\nu}_s(z) - \boldsymbol{\nu}_t(z)) = -\nabla_t \boldsymbol{\nu}(\mathbf{z}, t)$ be the partial derivative of the vector $\boldsymbol{\nu}(\mathbf{z}, t)$ w.r.t
 509 scalar t . Then we derive the diffusion loss, $\mathcal{L}_{\text{diffusion}}$, for the continuous case in the following manner

510 (for brevity we use the notation s for $s(i) = (i - 1)/T$ and t for $t(i) = i/T$):

$$\begin{aligned}
& \mathcal{L}_{\text{diffusion}} \\
&= \lim_{T \rightarrow \infty} \frac{1}{2} \sum_{i=2}^T \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} \text{D}_{\text{KL}}(q(\mathbf{x}_s | \mathbf{x}_t, \mathbf{x}_0, \mathbf{z}) \| p_{\theta}(\mathbf{x}_s | \mathbf{x}_t, \mathbf{z})) \\
&\quad \text{Using Eq. 64 we get,} \\
&= \lim_{T \rightarrow \infty} \frac{1}{2} \sum_{i=2}^T \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t(i)))^{\top} \text{diag}(\boldsymbol{\nu}(s(i), \mathbf{z}) - \boldsymbol{\nu}(t(i), \mathbf{z})) (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t(i))) \\
&= \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)} \left[\lim_{T \rightarrow \infty} \sum_{i=2}^T T (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t(i)))^{\top} \text{diag}(\boldsymbol{\nu}(s(i), \mathbf{z}) - \boldsymbol{\nu}(t(i), \mathbf{z})) (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t(i))) \frac{1}{T} \right] \\
&\quad \text{Using the fact that } \lim_{T \rightarrow \infty} T (\boldsymbol{\nu}(s, \mathbf{z}) - \boldsymbol{\nu}(z, t)) = -\nabla_t \boldsymbol{\nu}(t, \mathbf{z}) \text{ we get,} \\
&= -\frac{1}{2} \mathbb{E}_{t \sim \{0, \dots, 1\}} [(\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t))^{\top} (\nabla_t \boldsymbol{\nu}_t(z)) (\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t))] \\
&\quad \text{Substituting } \mathbf{x}_0 = \boldsymbol{\alpha}_t^{-1}(\mathbf{z})(\mathbf{x}_t - \boldsymbol{\sigma}_t(\mathbf{z})\boldsymbol{\epsilon}_t) \text{ from Eq. 53 and} \\
&\quad \text{Substituting } \mathbf{x}_{\theta}(\mathbf{x}_t, \mathbf{z}, t) = \boldsymbol{\alpha}_t^{-1}(\mathbf{z})(\mathbf{x}_t - \boldsymbol{\sigma}_t(\mathbf{z})\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)) \text{ from Eq. 54 we get,} \\
&= -\frac{1}{2} \mathbb{E}_{t \sim [0, 1]} \left[(\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t))^{\top} \left(\frac{\boldsymbol{\sigma}_t^2(\mathbf{z})}{\boldsymbol{\alpha}_t^2(\mathbf{z})} \times \nabla_t \boldsymbol{\nu}_t(\mathbf{z}) \right) (\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)) \right] \\
&\quad \text{Let } \boldsymbol{\nu}^{-1}(\mathbf{z}, t) \text{ denote the reciprocal of the values in the vector } \boldsymbol{\nu}(\mathbf{z}, t). \\
&= -\frac{1}{2} \mathbb{E}_{t \sim [0, 1]} [(\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t))^{\top} \text{diag}(\boldsymbol{\nu}^{-1}(t)(\mathbf{z}) \nabla_t \boldsymbol{\nu}_t(\mathbf{z})) (\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t))] \\
&\quad \text{Substituting } \boldsymbol{\nu}(\mathbf{z}, t) = \exp(-\boldsymbol{\gamma}(\mathbf{z}, t)) \text{ from Sec. D.1.1} \\
&= -\frac{1}{2} \mathbb{E}_{t \sim [0, 1]} [(\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t))^{\top} \text{diag}(\exp(\boldsymbol{\gamma}(\mathbf{z}, t)) \nabla_t \exp(-\boldsymbol{\gamma}(\mathbf{z}, t))) (\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t))] \\
&= \frac{1}{2} \mathbb{E}_{t \sim [0, 1]} [(\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t))^{\top} \text{diag}(\exp(\boldsymbol{\gamma}(\mathbf{z}, t)) \exp(-\boldsymbol{\gamma}(\mathbf{z}, t)) \nabla_t \boldsymbol{\gamma}(\mathbf{z}, t)) (\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t))] \\
&= \frac{1}{2} \mathbb{E}_{t \sim [0, 1]} [(\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t))^{\top} \text{diag}(\nabla_t \boldsymbol{\gamma}(\mathbf{z}, t)) (\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t))] \tag{65}
\end{aligned}$$

511 D.6 Recovering VDM from the Vectorized Representation of the diffusion loss

512 Notice that we recover the loss function in VDM when $\boldsymbol{\nu}(\mathbf{z}, t) = \nu(t)\mathbf{1}_{\mathbf{d}}$ where $\nu_t \in \mathbb{R}^+$ and $\mathbf{1}_{\mathbf{d}}$
513 represents a vector of 1s of size d and the noising schedule isn't conditioned on \mathbf{z} .

$$\begin{aligned}
& \int_0^1 \langle \mathbf{f}_{\theta}(\mathbf{x}_0, \boldsymbol{\nu}(\mathbf{z}, t)), \frac{d}{dt} \boldsymbol{\nu}(t) \rangle dt = \int_0^1 \langle \mathbf{f}_{\theta}(\mathbf{x}_0, \boldsymbol{\nu}(t)), \frac{d}{dt} (\nu(t)\mathbf{1}_{\mathbf{n}}) \rangle dt \\
&= \int_0^1 \langle \mathbf{f}_{\theta}(\mathbf{x}_0, \boldsymbol{\nu}(t)), \mathbf{1}_{\mathbf{d}} \rangle \nu'(t) dt \\
&= \int_0^1 \nu'(t) \|\mathbf{f}_{\theta}(\mathbf{x}_0, \boldsymbol{\nu}(t))\|_1 dt \\
&= \int_0^1 \nu'(t) \|(\mathbf{x}_0 - \tilde{\mathbf{x}}_{\theta}(\mathbf{x}_{\nu(t)}, \boldsymbol{\nu}(t)))\|_2^2 dt \tag{66}
\end{aligned}$$

514 $\int_0^1 \frac{d}{dt} \nu(t) \|(\mathbf{x}_0 - \tilde{\mathbf{x}}_{\theta}(\mathbf{x}_{\nu(t)}, \boldsymbol{\nu}(t)))\|_2^2 dt$ denotes the diffusion loss, $\mathcal{L}_{\text{diffusion}}$, as used in VDM; see
515 Kingma et al. (2021).

516 E Subset Sampling

517 Sampling a subset of k items from a collection of collection of n items, x_1, x_2, \dots, x_n belongs
518 to a category of algorithms called reservoir algorithms. In weighted reservoir sampling, every

519 x_i is associated with a weight $w_i \geq 0$. The probability associated with choosing the sequence
 520 $S_{\text{wrs}} = [i_1, i_2, \dots, i_k]$ be a tuple of indices. Then the probability associated with sampling this
 521 sequence is

$$p(S_{\text{wrs}}|\mathbf{w}) = \frac{w_{i_1}}{Z} \frac{w_{i_2}}{Z - w_{i_1}} \cdots \frac{w_{i_k}}{Z - \sum_{j=1}^{k-1} w_{i_j}} \quad (67)$$

522 [Efraimidis & Spirakis \(2006\)](#) give an algorithm for weighted reservoir sampling where each item
 523 is assigned a random key $r_i = u_i \frac{1}{w_i}$ where u_i is drawn from a uniform distribution $[0, 1]$ and w_i is
 524 the weight of item x_i . Let $\text{TopK}(\mathbf{r}, k)$ which takes keys $\mathbf{r} = [r_1, r_2, \dots, r_n]$ and returns a sequence
 525 $[i_1, i_2, \dots, i_k]$. [Efraimidis & Spirakis \(2006\)](#) proved that $\text{TopK}(\mathbf{r}, k)$ is distributed according to
 526 $p(S_{\text{wrs}}|\mathbf{w})$.

527 Let’s represent a subset $S \in \{0, 1\}^n$ with exactly k non-zero elements that are equal to 1. Then the
 528 probability associated with sampling S is given as,

$$p(S|\mathbf{w}) = \sum_{S_{\text{wrs}} \in \Pi(S)} p(S_{\text{wrs}}|\mathbf{w}) \quad (68)$$

529 where $\Pi(S)$ denotes all possible permutations of the sequence S . By ignoring the ordering of the ele-
 530 ments in S_{wrs} we can sample using the same algorithm. [Xie & Ermon \(2019\)](#) show that this sampling
 531 algorithm is equivalent to $\text{TopK}(\hat{\mathbf{r}}, k)$ where $\hat{\mathbf{r}} = [\hat{r}_1, \hat{r}_2, \dots, \hat{r}_n]$ where $\hat{r}_i = -\log(-\log(r_i)) =$
 532 $\log w_i + \text{Gumbel}(0, 1)$. This holds true because the monotonic transformation $-\log(-\log(x))$
 533 preserves the ordering of the keys and thus $\text{TopK}(\mathbf{r}, k) \equiv \text{TopK}(\hat{\mathbf{r}}, k)$.

534 **Sum of Gamma Distribution.** [Niepert et al. \(2021\)](#) show that adding SOG noise instead of Gumbel
 535 noise leads to better performance.

536 [Niepert et al. \(2021\)](#) show that $\mathbf{z} \sim p_\theta(\mathbf{z}; \theta)$ is equivalent to $\mathbf{z} = \arg \max_{y \in Y} \langle \theta + \epsilon_g, y \rangle$ where ϵ_g is
 537 a sample from Sum-of-Gamma distribution given by

$$\text{SoG}(k, \tau, s) = \frac{\tau}{k} \left(\sum_{i=1}^s \text{Gamma}\left(\frac{1}{k}, \frac{k}{i}\right) - \log s \right), \quad (69)$$

538 where s is a positive integer and $\text{Gamma}(\alpha, \beta)$ is the Gamma distribution with (α, β) as the shape
 539 and scale parameters.

540 And hence, given logits $\log \mathbf{w}$, we sample a k -hot vector using $\text{TopK}(\log \mathbf{w} + \epsilon)$. We choose a
 541 categorical prior with uniform distribution across n classes. Thus the KL loss term is given by:

$$-\sum_{i=1}^n \frac{w_i}{Z} \log \left(n \frac{w_i}{Z} \right) \quad (70)$$

542 F Experiment Details

543 F.1 Model Architecture

544 **Denosing network.** Our model architecture is extremely similar to VDM. The UNet of our pixel-
 545 space diffusion has an unchanged architecture from [Kingma et al. \(2021\)](#). This structure is specifically
 546 designed for optimal performance in maximum likelihood training. We employ features from VDM
 547 such as the elimination of internal downsampling/upsampling processes and the integration of Fourier
 548 features to enhance fine-scale prediction accuracy. In alignment with the configurations suggested by
 549 [Kingma et al. \(2021\)](#), our approach varies depending on the dataset: For CIFAR-10, we employ a
 550 U-Net with a depth of 32 and 128 channels; for ImageNet-32, the U-Net also has a depth of 32, but
 551 the channel count is increased to 256. Additionally, all these models incorporate a dropout rate of 0.1
 552 in their intermediate layers.

553 **Encoder network.** $q_\phi(\mathbf{z}|\mathbf{x})$ is modeled using a sequence of 4 Resnet blocks with a channel count
 554 of 128 for CIFAR-10 and 256 for ImageNet-32 with a drop out of 0.1 in their intermediate layers.

Table 3: Likelihood in bits per dimension (BPD) on the test set of CIFAR-10 and ImageNet. Results with “/” means they are not reported in the original papers. Model types are autoregressive (AR), normalizing flows (Flow), variational autoencoders (VAE), diffusion models (Diff), diffusion ODEs (Diff ODE). The likelihood for “Diff” type models is computed using the VLB-based method described in appendix Sec. H.1, while “Diff ODE” type models utilize an ODE-based exact likelihood estimate as detailed in appendix Sec. H.2. Additionally, for MULAN, we present the mean and a 95% confidence interval.

Model	Type	CIFAR-10 (\downarrow)	ImageNet (\downarrow)
PixelCNN (Van den Oord et al., 2016)	AR	3.03	3.83
PixelCNN++ (Salimans et al., 2017)	AR	2.92	/
Glow (Kingma & Dhariwal, 2018)	Flow	/	4.09
Image Transformer (Parmar et al., 2018)	AR	2.90	3.77
DDPM (Ho et al., 2020)	Diff	3.69	/
Score SDE (Song et al., 2020)	Diff	2.99	/
Improved DDPM (Nichol & Dhariwal, 2021)	Diff	2.94	/
VDM (Kingma et al., 2021)	Diff	2.65	3.72
Flow Matching (Lipman et al., 2022)	Flow	2.99	/
i-DODE (Zheng et al., 2023) (VLB-based)	Diff	2.61	/
i-DODE (Zheng et al., 2023) (ODE-based)	Diff ODE	2.56	3.69
MULAN (Ours, VLB-based; see Sec. H.1)	Diff	2.60	3.71
MULAN (Ours, ODE-based; see Sec. H.2)	Diff ODE	2.55 ± 10^{-3}	3.67 ± 10^{-3}

555 **Noise schedule.** For polynomial noise schedule, we use an MLP that maps the latent vector \mathbf{z}
 556 to $\mathbf{a}_\phi(\mathbf{z})$, $\mathbf{b}_\phi(\mathbf{z})$, $\mathbf{c}(\mathbf{z})$; see Eq. D.2 for details. The MLP has 2 hidden layers of size 3072 with
 557 `swish` activation function. The final layer is a linear mapping to 3×3072 values corresponding to
 558 $\mathbf{a}_\phi(\mathbf{z})$, $\mathbf{b}_\phi(\mathbf{z})$, $\mathbf{c}(\mathbf{z})$. Note that $\mathbf{a}_\phi(\mathbf{z})$, $\mathbf{b}_\phi(\mathbf{z})$, $\mathbf{c}(\mathbf{z})$ have the same dimensionality of 3072.

559 F.2 Hardware.

560 For the ImageNet experiments, we used a single GPU node with 8-A100s. For the cifar-10 experi-
 561 ments, we used several GPUs types including V100, A5000s, A6000s, A100-40GBs, and 3090s but
 562 the experiments were trained on 4 GPUs with `float32` precision.

563 F.3 Training

564 This section reports experimental results on the CIFAR-10 (Krizhevsky et al., 2009) and ImageNet-
 565 32 (Van Den Oord et al., 2016) datasets. We chose to employ a discrete prior for the auxiliary latent
 566 space rather than a Gaussian prior due to training instability issues that frequently led to NaNs. In all
 567 our experiments, we set the parameters for the discrete latent distribution as $m = 50$ and $k = 15$.

568 F.4 Likelihood Estimation.

569 In Table 1, we present the likelihood estimation results for MULAN, and other recent methods on
 570 CIFAR-10 and ImageNet-32 using the VLB-estimate; details in Sec. H.1. This version of MULAN
 571 was trained with noise parameterization for 10M steps on CIFAR-10 and 2M steps on Imagenet-32,
 572 similar to VDM (Kingma et al., 2021). We apply MULAN on top of the VDM, endowing it with a
 573 learned multivariate noising schedule conditioned on auxiliary latent variables. We find that these
 574 new components result in a significant improvement in BPD over a vanilla VDM.

575 We also compute the likelihood using ODE-based exact likelihood estimate with which we outperform
 576 all existing methods in density estimation on CIFAR-10 and ImageNet-32. We train MULAN for 8M
 577 steps using velocity parameterization on CIFAR-10 and 2M steps on Imagenet-32. During inference,
 578 we extract the underlying probability flow ODE; see Sec. H.2, just like Zheng et al. (2023). Although
 579 Zheng et al. (2023) used various other techniques, such as importance-sampled training, variance
 580 minimization, and higher-order score matching, MULAN did not incorporate these. Combining these
 581 strategies with MULAN could further improve its performance.

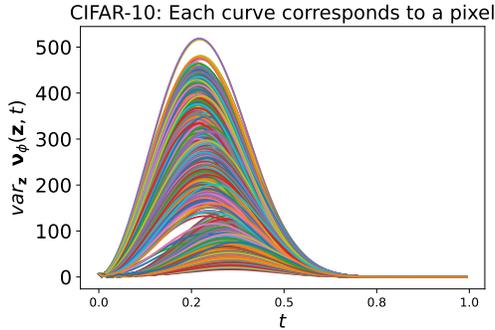
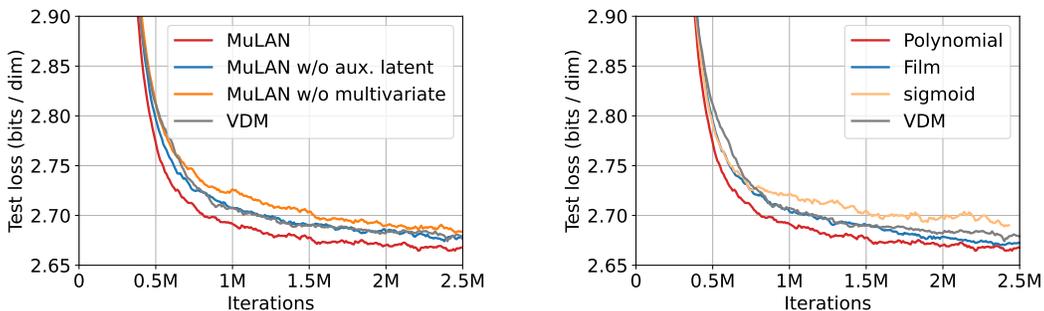


Figure 2: Noise schedule visualizations for MULAN on CIFAR-10. In this figure, we plot the variance of $\nu_\phi(\mathbf{z}, t)$ across different $\mathbf{z} \sim p_\theta(\mathbf{z})$ where each curve represents the SNR corresponding to an input dimension.



(a) Ablations for different components of MULAN. MULAN w/o aux. latent refers to MULAN where the noise schedule isn't conditioned on an auxiliary latent variable. MULAN w/o multivariate refers to MULAN with a scalar noise schedule. VDM uses a scalar noise schedule and doesn't have an auxiliary latent space. We see that MULAN performs the best.

(b) Ablations for noise schedules with different functional parameterizations: sigmoid, polynomial, and a monotonic neural network with FiLM conditioning. We observe that the polynomial parameterization performs the best.

Figure 3: CIFAR-10 ablation studies with a reduced batch size and fewer training steps.

582 F.5 Ablation Analysis

583 Due to the expensive cost of training, we only performed ablation studies on CIFAR-10 with a reduced
 584 batch size of 64 and trained the model for 2.5M training steps. In Fig. 3 we ablate each component of
 585 MULAN: when we remove the conditioning on an auxiliary latent space from MULAN so that we
 586 have a multivariate noise schedule that is solely conditioned on time t , our performance becomes
 587 comparable to that of VDM, on which our model is based. Modifying our method to have a scalar
 588 noise schedule conditioned on the auxiliary latent variable \mathbf{z} leads to slightly lower performance than
 589 VDM in the initial training stages. However, it gradually converges toward VDM.

590 **Loss curves for different noise schedules.** We investigate different parameterizations of the noise
 591 schedule in Fig. 3. Among polynomial, sigmoid, and monotonic neural network, we find that the
 592 polynomial parameterization yields the best performance. The polynomial noise schedule is a novel
 593 component introduced in our work.

594 **Replacing the noise schedules in a trained denoising model.** We also wish to confirm experimen-
 595 tally our claim that the learning objective is not invariant to our choice of multivariate noise schedule.
 596 To investigate this, we replace the noise schedule in the trained denoising model with two alternatives:
 597 MULAN with scalar noise schedule, and a linear noise schedule: $\gamma_\phi(\mathbf{z}, t) = \gamma_{\min} + t(\gamma_{\max} - \gamma_{\min})\mathbf{1}_d$;
 598 see (Kingma et al., 2021). For both the noise schedules the likelihood worsens to the same value as
 599 that of the VDM: 2.65. This experimental result strongly supports our theory that all scalar noise
 600 schedules are equivalent, as they compute the likelihood along the same diffusion trajectory. It also

601 underscores that it is not the multivariate nature or the auxiliary latent space individually, but the
602 combination of both, that makes MULAN effective.

603 **Examining the noise schedule.** Since the noise schedule, $\gamma_\phi(\mathbf{z}, t)$ is multivariate, we expect to
604 learn different noise schedules for different input dimensions and different inputs $\mathbf{z} \sim p_\theta(\mathbf{z})$. In Fig. 2,
605 we take our best trained model on CIFAR-10 and visualize the variance of the noise schedule at each
606 point in time for different pixels, where the variance is taken on 128 samples $\mathbf{z} \sim p_\theta(\mathbf{z})$. We note
607 an increased variation in the early portions of the noise schedule. However, on an absolute scale,
608 the variance of this noise is smaller than we expected. We also tried to visualize noise schedules
609 across different dataset images and across different areas of the same image; refer to Fig. 11. We
610 also generated synthetic datasets in which each datapoint contained only high frequencies or only low
611 frequencies, and with random masking applied to parts of the data points; see Sec. G. Surprisingly,
612 none of these experiments revealed human-interpretable patterns in the learned schedule, although
613 we did observe clear differences in likelihood estimation. We hypothesize that other architectures and
614 other forms of conditioning may reveal interpretable patterns of variation; however, we leave this
615 exploration to future work.

616 G Datasets and Visualizations

617 In this section we provide a brief description of the datasets used in the paper and visualize the
618 generated samples and the noise schedules.

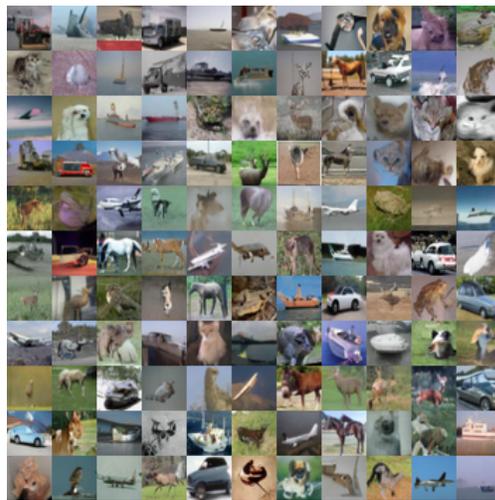
619 G.1 CIFAR-10

620 The CIFAR-10 dataset (Krizhevsky et al., 2009) is a collection of images consisting of 60,000 32×32
621 color images in 10 different classes, with each class representing a distinct object or scene. The
622 dataset is divided into 50,000 training images and 10,000 test images, with each class having an equal
623 representation in both sets. The classes in CIFAR-10 include: Airplane, Automobile, Bird, Cat, Deer,
624 Dog, Frog, Horse, Ship, Truck.

625 Randomly generated samples for the CIFAR-10 dataset are provided in Fig. 4a for MULAN and
626 Fig. 4b for VDM. We visualize the noise schedule in Fig. 11.



(a) MULAN with velocity reparameterization after 8M training iterations.



(b) VDM after 10M training iterations.

Figure 4: CIFAR-10 samples generated by different methods.

627 **G.2 ImageNet-32**

628 ImageNet-32 is a dataset derived from ImageNet [Deng et al. \(2009\)](#), where the original images have
629 been resized to a resolution of 32×32 . This dataset comprises 1,281,167 training samples and 50,000
630 test samples, distributed across 1,000 labels.

631 Randomly generated samples for the ImageNet dataset are provided in [Fig. 5](#) for MULAN and
632 [Fig. 6](#) for VDM. We visualize the noise schedule in [Fig. 11](#).

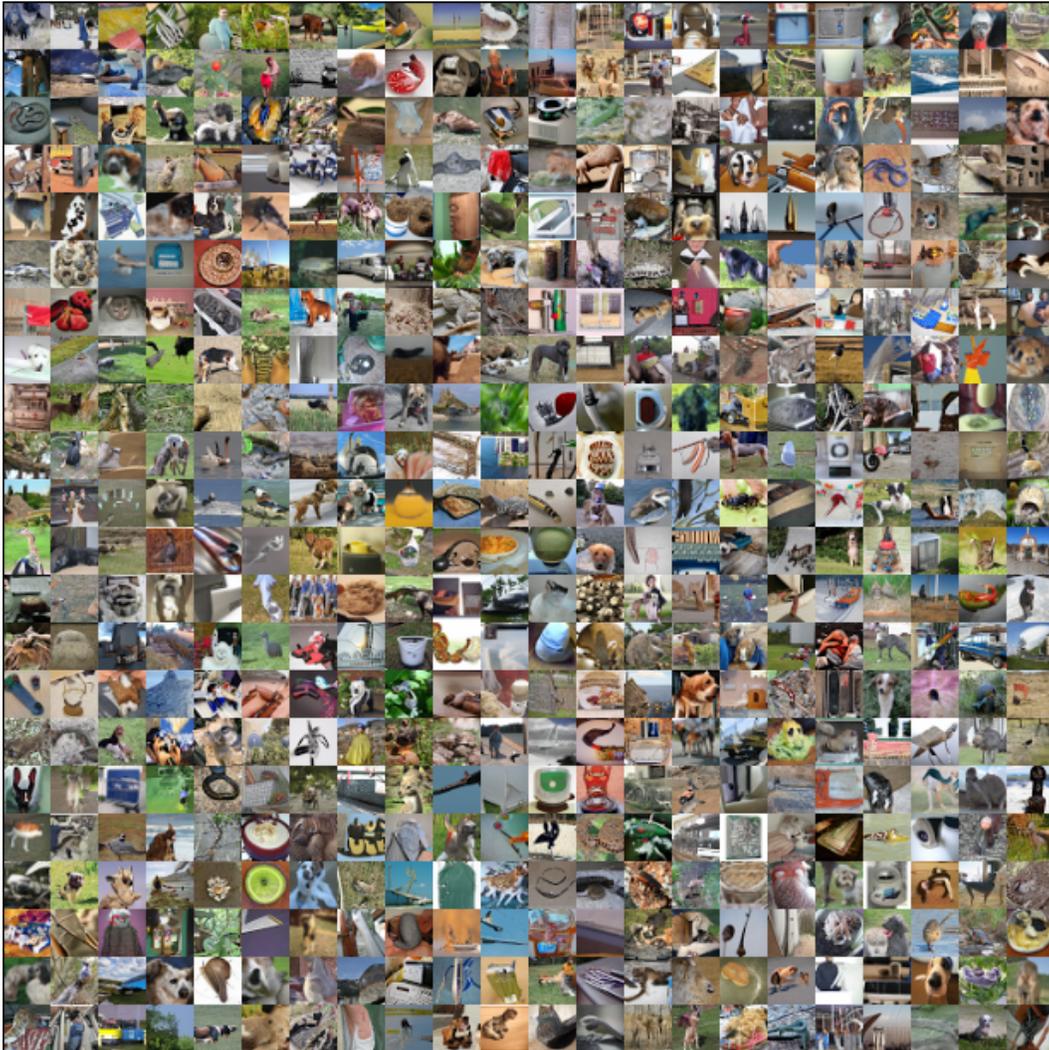


Figure 5: MULAN with noise parameterization after 2M training iterations.

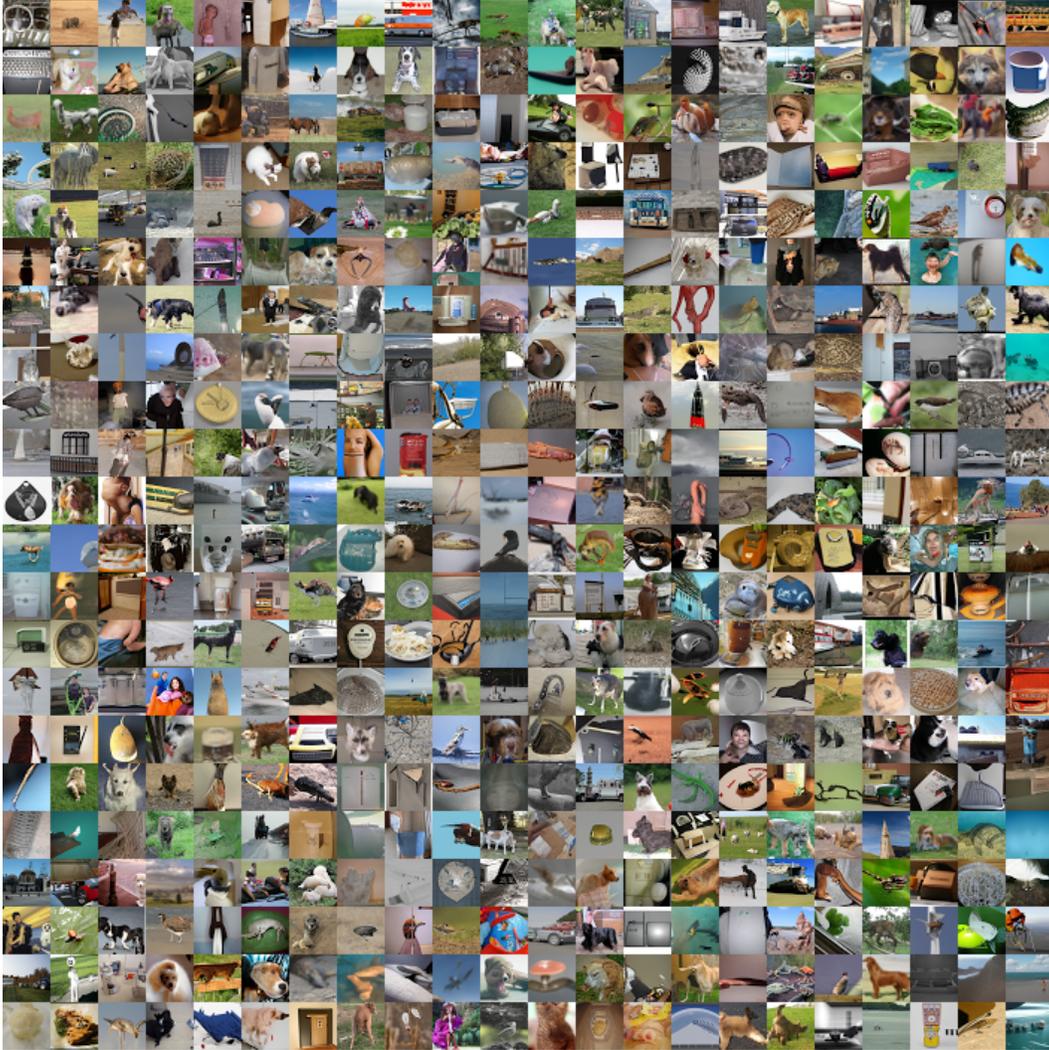


Figure 6: VDM after 2M training iterations.

633 **G.3 Frequency**

634 To see if MULAN learns different noise schedules for images with different intensities, we modify
 635 the images in the CIFAR-10 dataset where we modify an image where we randomly remove the low
 636 frequency component an image or remove the high frequency with equal probability. Fig. 7a shows
 637 the training samples. MULAN was trained for 500K steps. The samples generated by MULAN is
 638 shown in Fig. 7b. The corresponding noise schedules is shown in Fig. 11. As compared to CIFAR-10,
 639 we notice that the spatial variation in the noise schedule increases (SNRs for all the pixels form a
 640 wider band) while the variance of the noise schedule across instances decreases slightly.



(a) Training samples.

(b) Samples generated by MULAN with noise parameterization after 500K training iterations.

Figure 7: Frequency Split CIFAR-10 dataset.

641 G.4 Frequency-2

642 To see if MULAN learns different noise schedules for images with different intensities, we modify
 643 the images in the CIFAR-10 dataset where we modify an image where we randomly remove the low
 644 frequency component an image or remove the high frequency with equal probability. Fig. 7a shows
 645 the training samples. MULAN was trained for 500K steps. The samples generated by MULAN is
 646 shown in Fig. 7b. The corresponding noise schedules is shown in Fig. 11. As compared to CIFAR-10,
 647 we notice that the spatial variation in the noise schedule increases (SNRs for all the pixels form a
 648 wider band) and the variance of the noise schedule across instances increases as well.



(a) Training samples.

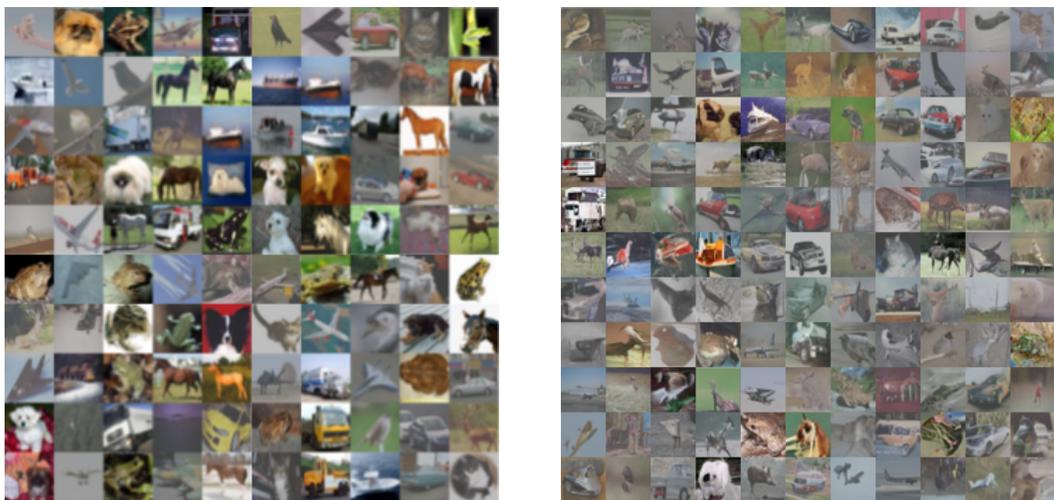
(b) Samples generated by MULAN with noise parameterization after 500K training iterations.

Figure 8: Frequency Split-2 CIFAR-10 dataset.

649 G.5 CIFAR-10: Intensity

650 To see if MULAN learns different noise schedules for images with different intensities, we modify
 651 the images in the CIFAR-10 dataset where we randomly convert an image into a low intensity or

652 a high intensity image with equal probability. Originally, the CIFAR10 images are in the range [0,
 653 255]. To convert an image into a low intensity image we multiply all pixel values by 0.5. To convert
 654 an image into a high intensity image we multiply all the pixel values by 0.5 and add 127.5 to them.
 655 Fig. 9a shows the training samples. MULAN was trained for 500K steps. The samples generated by
 656 MULAN is shown in Fig. 9b. The corresponding noise schedules is shown in Fig. 11. As compared
 657 to CIFAR-10, we notice that the spatial variation in the noise schedule slightly increases (SNRs for
 658 all the pixels form a wider band) while the variance of the noise schedule across instances slightly
 659 decreases.



(a) Training samples.

(b) Samples generated by MULAN with noise parameterization after 500K training iterations.

Figure 9: Intensity CIFAR-10 dataset.

660 G.6 Mask

661 We modify the CIFAR-10 dataset where we randomly mask (i.e. replace with 0s) the top of an
 662 image or the bottom half of an image with equal probability. Fig. 10a shows the training samples.
 663 MULAN was trained for 500K steps. The samples generated by MULAN is shown in Fig. 10b. The
 664 corresponding noise schedules is shown in Fig. 11. As compared to CIFAR-10, we notice that the
 665 spatial variation in the noise schedule slightly increases (SNRs for all the pixels form a wider band)
 666 while the variance of the noise schedule across instances decreases.



(a) Training samples.



(b) Samples generated by MULAN with noise parameterization after 500K training iterations.

Figure 10: Intensity CIFAR-10 dataset.

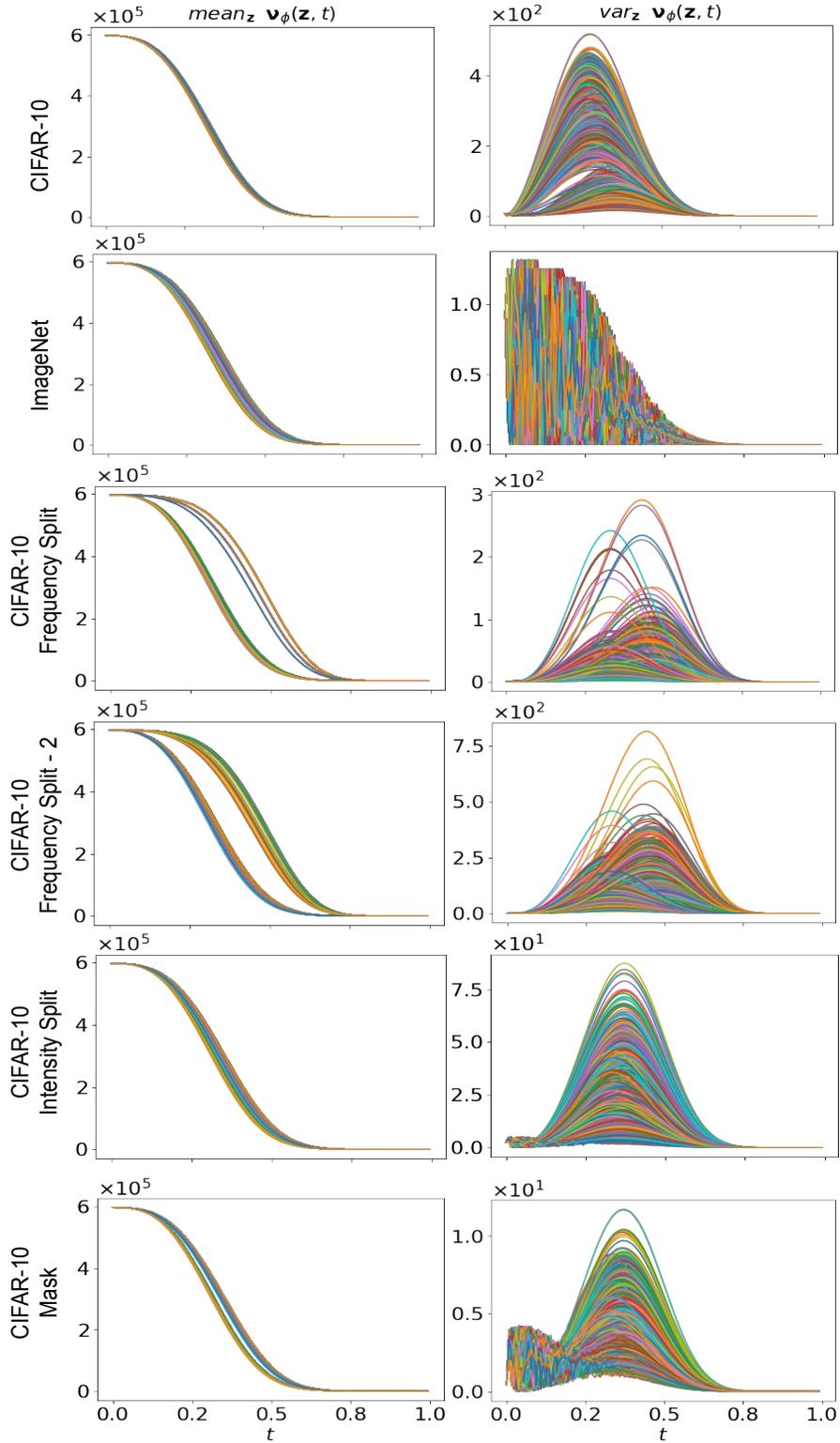


Figure 11: signal-to-noise ratio for different datasets.

667 H Likelihood Estimation

668 We used both Variance Lower Bound (VLB) and ODE-based methods to compute BPD.

669 H.1 VLB Estimate

670 In the VLB-based approach, we employ Eq. 60. To compute $\mathcal{L}_{\text{diffusion}}$, we use $T = 128$ in Eq. 61,
671 discretizing the timesteps, $t \in [0, 1]$ into 128 bins.

672 H.2 Exact likelihood computation using Probability Flow ODE

673 A diffusion process whose marginal is given by (the same as in Eq. 71),

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\alpha}_t\mathbf{x}_0, \text{diag}(\boldsymbol{\sigma}_t^2)), \quad (71)$$

674 can be modeled as the solution to an Itô Stochastic Differential Equation (SDE):

$$d\mathbf{x}_t = \mathbf{f}(t)\mathbf{x}_t dt + \mathbf{g}(t)d\mathbf{w}_t, \quad \mathbf{x}_0 \sim q_0(\mathbf{x}_0), \quad (72)$$

675 where $\mathbf{f}(t) \in \mathbb{R}^d$, $\mathbf{g}(t) \in \mathbb{R}^d$ take the following expressions (Song et al., 2020):

$$\begin{aligned} \mathbf{f}(t) &= \frac{d}{dt} \log \boldsymbol{\alpha}_t, \\ \mathbf{g}^2(t) &= \frac{d}{dt} \boldsymbol{\sigma}_t^2 - 2\boldsymbol{\sigma}_t^2 \frac{d}{dt} \log \boldsymbol{\alpha}_t \end{aligned}$$

676 The corresponding reverse process, Eq. 2, can also be modelled by an equivalent reverse-time SDE:

$$d\mathbf{x}_t = [\mathbf{f}(t) - \mathbf{g}(t)^2 \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0)]dt + \mathbf{g}(t)d\bar{\mathbf{w}}_t, \quad \mathbf{x}_1 \sim p_\theta(\mathbf{x}_1), \quad (73)$$

677 where $\bar{\mathbf{w}}$ is a standard Wiener process when time flows backwards from $1 \rightarrow 0$, and dt is an
678 infinitesimal negative timestep. Song et al. (2020) show that the marginals of Eq. 73 can be described
679 by the following Ordinary Differential Equation (ODE) in the reverse process:

$$d\mathbf{x}_t = \left[\mathbf{f}(t)\mathbf{x}_t - \frac{1}{2}\mathbf{g}^2(t)\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0) \right] dt. \quad (74)$$

680 This ODE, also called the probability flow ODE, allows us to compute the exact likelihood on any
681 input data via the instantaneous change of variables formula as proposed in Chen et al. (2018). Note
682 that during the reverse process, the term $q(\mathbf{x}_t|\mathbf{x}_0)$ is unknown and is approximated by parameterized
683 by $p_\theta(\mathbf{x}_t)$. For the probability flow defined in Eq. 74, Chen et al. (2018) show that the log-likelihood
684 of $p_\theta(\mathbf{x}_0)$ can be computed using the following equation:

$$\log p_\theta(\mathbf{x}_0) = \log p_\theta(\mathbf{x}_1) - \int_{t=0}^{t=1} \text{tr}(\nabla_{\mathbf{x}_t} \mathbf{h}_\theta(\mathbf{x}_t, t)) dt, \quad (75)$$

$$\text{where } \mathbf{h}_\theta(\mathbf{x}_t, t) \equiv \mathbf{f}(t)\mathbf{x}_t - \frac{1}{2}\mathbf{g}^2(t)\nabla_{\mathbf{x}_t} \log p_\theta(\mathbf{x}_t)$$

685 H.2.1 Probability Flow ODE for MULAN.

686 Similarly for the forward process conditioned on the auxiliary latent variable, \mathbf{z} ,

$$q_\phi(\mathbf{x}_t|\mathbf{x}_0, \mathbf{z}) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\alpha}_t(\mathbf{z})\mathbf{x}_0, \text{diag}(\boldsymbol{\sigma}_t^2(\mathbf{z}))), \quad \mathbf{x}_0 \sim q_0(\mathbf{x}_0), \quad \mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}_0), \quad (76)$$

687 we can extend Eq. 72 in the following manner,

$$d\mathbf{x}_t = \mathbf{f}(\mathbf{z}, t)\mathbf{x}_t dt + \mathbf{g}(\mathbf{z}, t)d\mathbf{w}_t, \quad \mathbf{x}_0 \sim q_0(\mathbf{x}_0), \quad \mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}_0), \quad (77)$$

688 to obtain the corresponding SDE formulation. Notice that the random variable \mathbf{z} in the above equation
689 doesn't have a subscript t , and hence, \mathbf{z} is drawn from $q_\phi(\mathbf{z}|\mathbf{x}_0)$ once and the same \mathbf{z} is used as \mathbf{x}_0
690 diffuses to \mathbf{x}_1 . The expressions for $\mathbf{f}(\mathbf{z}, t) : \mathbb{R}^m \times [0, 1] \rightarrow \mathbb{R}^d$, $\mathbf{g}(\mathbf{z}, t) : \mathbb{R}^m \times [0, 1] \rightarrow \mathbb{R}^d$ is given
691 as follows:

$$\begin{aligned} \mathbf{f}(\mathbf{z}, t) &= \frac{d}{dt} \log \boldsymbol{\alpha}_t(\mathbf{z}), \\ \mathbf{g}^2(\mathbf{z}, t) &= \frac{d}{dt} \boldsymbol{\sigma}_t^2(\mathbf{z}) - 2\boldsymbol{\sigma}_t^2(\mathbf{z}) \frac{d}{dt} \log \boldsymbol{\alpha}_t(\mathbf{z}) \end{aligned}$$

692 Recall that $\alpha_t^2(\mathbf{z}) = \text{sigmoid}(-\gamma_\phi(\mathbf{z}, t))$, $\sigma_t^2(\mathbf{z}) = \text{sigmoid}(\gamma_\phi(\mathbf{z}, t))$. Substituting these in the
 693 above equations, the expressions for $\mathbf{f}(\mathbf{z}, t)$ and $\mathbf{g}^2(\mathbf{z}, t)$ simplify to the following:

$$\begin{aligned}\mathbf{f}(\mathbf{z}, t) &= -\frac{1}{2}\text{sigmoid}(\gamma_\phi(\mathbf{z}, t))\frac{d}{dt}\gamma_\phi(\mathbf{z}, t), \\ \mathbf{g}^2(\mathbf{z}, t) &= \text{sigmoid}(\gamma_\phi(\mathbf{z}, t))\frac{d}{dt}\gamma_\phi(\mathbf{z}, t)\end{aligned}$$

694 The corresponding reverse-time SDE is given as:

$$d\mathbf{x}_t = [\mathbf{f}(t) - \mathbf{g}(t)^2 \nabla_{\mathbf{x}_t} \log q_\phi(\mathbf{x}_t | \mathbf{x}_0, \mathbf{z})] dt + \mathbf{g}(t) d\bar{\mathbf{w}}_t, \quad \mathbf{x}_1 \sim p_\theta(\mathbf{x}_1), \quad \mathbf{z} \sim p_\theta(\mathbf{z}), \quad (78)$$

695 where $\bar{\mathbf{w}}$ is a standard Wiener process when time flows backwards from $1 \rightarrow 0$, and dt is an
 696 infinitesimal negative timestep. Given, $\mathbf{s}_\theta(\mathbf{x}_t, \mathbf{z})$, an approximation to the true score function,
 697 $\nabla_{\mathbf{x}_t} \log q_\phi(\mathbf{x}_t | \mathbf{x}_0, \mathbf{z})$, Song et al. (2020) show that the marginals of Eq. 78 can be described by the
 698 following Ordinary Differential Equation (ODE):

$$d\mathbf{x}_t = \left[\mathbf{f}(\mathbf{z}, t) - \frac{1}{2} \mathbf{g}^2(\mathbf{z}, t) \mathbf{s}_\theta(\mathbf{x}_t, \mathbf{z}) \right] dt, \quad (79)$$

Zheng et al. (2023) show that the score function, $\mathbf{s}_\theta(\mathbf{x}_t, \mathbf{z})$, for the noise and the velocity parameteri-
 zation is given as follows:

$$\mathbf{s}_\theta(\mathbf{x}_t, \mathbf{z}) = \begin{cases} -\frac{\epsilon_\theta(\mathbf{x}_t, t)}{\sigma_t(\mathbf{z})} & \text{Noise parameterization; see Sec. D.1.1} \quad (80a) \\ -\mathbf{x}_t - \exp\left(-\frac{1}{2}\gamma_\phi(\mathbf{z}, t)\right) \mathbf{v}_\theta(\mathbf{x}_t, \mathbf{z}, t) & \text{Velocity parameterization; see Sec. D.1.2} \quad (80b) \end{cases}$$

699 Applying the change of variables formula (Chen et al., 2018) on Eq. 79, $\log p_\theta(\mathbf{x}_0 | \mathbf{z})$ can be computed
 700 in the following manner:

$$\begin{aligned}\log p_\theta(\mathbf{x}_0 | \mathbf{z}) &= \log p_\theta(\mathbf{x}_1) - \int_{t=0}^{t=1} \text{tr}(\nabla_{\mathbf{x}_t} \mathbf{h}_\theta(\mathbf{x}_t, \mathbf{z}, t)) dt, \quad (81) \\ \text{where } \mathbf{h}_\theta(\mathbf{x}_t, \mathbf{z}, t) &\equiv \mathbf{f}(\mathbf{z}, t) - \frac{1}{2} \mathbf{g}^2(\mathbf{z}, t) \mathbf{s}_\theta(\mathbf{x}_t, \mathbf{z})\end{aligned}$$

701 The expression for log-likelihood (Eq. 6) is as follows,

$$\begin{aligned}\log p_\theta(\mathbf{x}_0) &\geq \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0 | \mathbf{z})] - \text{D}_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}_0) \| p_\theta(\mathbf{z})) \\ \text{Using Eq. 81,} \\ &= \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x}_0)} \left[\log p_\theta(\mathbf{x}_1) - \int_{t=0}^{t=1} \text{tr}(\nabla_{\mathbf{x}_t} \mathbf{h}_\theta(\mathbf{x}_t, t, \mathbf{z})) dt \right] - \text{D}_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}_0) \| p_\theta(\mathbf{z}))\end{aligned} \quad (82)$$

702 Computing $\text{tr}(\nabla_{\mathbf{x}_t} \mathbf{h}_\theta(\mathbf{x}_t, t, \mathbf{z}))$ is expensive and we follow Chen et al. (2018); Zheng et al. (2023);
 703 Grathwohl et al. (2018) to estimate it with Skilling-Hutchinson trace estimator (Skilling, 1989;
 704 Hutchinson, 1989). In particular, we have

$$\text{tr}(\nabla_{\mathbf{x}_t} \mathbf{h}_\theta(\mathbf{x}_t, t, \mathbf{z})) = \mathbb{E}_{p(\epsilon)} [\epsilon^\top \nabla_{\mathbf{x}_t} \mathbf{h}_\theta(\mathbf{x}_t, t, \mathbf{z}) \epsilon], \quad (83)$$

705 where the random variable ϵ satisfies $\mathbb{E}_{p(\epsilon)}[\epsilon] = \mathbf{0}$ and $\text{Cov}_{p(\epsilon)}[\epsilon] = \mathbf{I}$. Common choices for $p(\epsilon)$
 706 include Rademacher or Gaussian distributions. Notably, the term $\nabla_{\mathbf{x}_t} \mathbf{h}_\theta(\mathbf{x}_t, t, \mathbf{z}) \epsilon$ can be computed
 707 efficiently using ‘‘Jacobian-vector-product’’ computation in JAX. In our experiments, we follow the
 708 exact evaluation procedure for computing likelihood as outlined in Song et al. (2020); Grathwohl
 709 et al. (2018). Specifically, for the computation of Eq. 83, we employ a Rademacher distribution for
 710 $p(\epsilon)$. To calculate the integral in Eq. 82, we utilize the RK45 ODE solver (Dormand & Prince, 1980)
 711 provided by `scipy.integrate.solve_ivp` with `atol=1e-5` and `rtol=1e-5`.

712 **H.2.2 Dequantization.**

Real-world datasets for images or texts often consist of discrete data. Attempting to learn a continuous density model directly on these discrete data points can lead to degenerate outcomes (Uria et al., 2013) and fail to provide meaningful density estimations. Dequantization (Salimans et al., 2017; Ho et al., 2020; Zheng et al., 2023) is a common solution in such cases. To elaborate, let x_0 represent 8-bit discrete data scaled to $[-1, 1]$. Dequantization methods assume that we have trained a continuous model distribution p_θ for x_0 , and define the discrete model distribution by

$$P_\theta(\mathbf{x}_0) = \int_{[-\frac{1}{256}, \frac{1}{256}]^d} p_\theta(\mathbf{x}_0 + u) du.$$

713 To train $P_\theta(\mathbf{x}_0)$ by maximum likelihood estimation, variational dequantization (Ho et al., 2020;
714 Zheng et al., 2023) introduces a dequantization distribution $q(u|\mathbf{x}_0)$ and jointly train p_{model} and
715 $q(u|\mathbf{x}_0)$ by a variational lower bound:

$$\log P_\theta(\mathbf{x}_0) \geq \mathbb{E}_{q(u|\mathbf{x}_0)} [p_\theta(\mathbf{x}_0 + u) - \log q(u|\mathbf{x}_0)]. \quad (84)$$

Truncated Normal Dequantization. Zheng et al. (2023) show that truncated Normal distribution,

$$q(u|\mathbf{x}_0) = \mathcal{TN}\left(\mathbf{0}, \mathbf{I}, -\frac{1}{256}, \frac{1}{256}\right)$$

716 with mean $\mathbf{0}$, covariance \mathbf{I} , and bounds $[-\frac{1}{256}, \frac{1}{256}]$ along each dimension, leads to a better likelihood
717 estimate. Thus, Eq. 84 simplifies to the following (for details please refer to section A. in Zheng et al.
718 (2023)):

$$\log P_\theta(\mathbf{x}_0) \geq \mathbb{E}_{\hat{\epsilon} \sim \mathcal{TN}(\mathbf{0}, \mathbf{I}, -\tau, \tau)} \left[\log p_\theta\left(\mathbf{x}_0 + \frac{\sigma_\epsilon}{\alpha_\epsilon} \hat{\epsilon}\right) \right] + \frac{d}{2} (1 + \log(2\pi\sigma_\epsilon^2)) - 0.01522 \times d \quad (85)$$

$$\text{with } \frac{\sigma_\epsilon}{\alpha_\epsilon} = \exp\left(-\frac{1}{2} \times 13.3\right),$$

$$\sigma_\epsilon = \text{sqrt}(\text{sigmoid}(-13.3)), \text{ and } \tau = 3.$$

719 $\log p_\theta\left(\mathbf{x}_0 + \frac{\sigma_\epsilon}{\alpha_\epsilon} \hat{\epsilon}\right)$ is evaluated using Eq. 82.

720 **Importance Weighted Estimator.** Eq. 85 can also be extended to obtain an importance weighted
721 likelihood estimator to get a tighter bound on the likelihood. The variational bound is given by (for
722 details please refer to section A. in Zheng et al. (2023)):

$$\log P_\theta(\mathbf{x}_0) \geq \mathbb{E}_{\hat{\epsilon}^{(1)}, \dots, \hat{\epsilon}^{(K)} \sim \mathcal{TN}(\mathbf{0}, \mathbf{I}, -\tau, \tau)} \left[\log \left(\frac{1}{K} \sum_{i=1}^K \frac{p_\theta\left(\mathbf{x}_0 + \frac{\sigma_\epsilon}{\alpha_\epsilon} \hat{\epsilon}^{(i)}\right)}{q(\hat{\epsilon}^{(i)})} \right) \right] + d \log \sigma_\epsilon \quad (86)$$

$$\text{with } \frac{\sigma_\epsilon}{\alpha_\epsilon} = \exp\left(-\frac{1}{2} \times 13.3\right), \log \sigma_\epsilon = \frac{1}{2}(-13.3 + \text{softplus}(-13.3)),$$

$$q(\hat{\epsilon}) = \frac{1}{(2\pi Z)^2} \exp\left(-\frac{1}{2} \|\hat{\epsilon}\|_2^2\right), Z = 0.9974613, \text{ and } \tau = 3.$$

723 Note that for $K = 1$, Eq. 86 is equivalent to Eq. 85; see Zheng et al. (2023). $\log p_\theta\left(\mathbf{x}_0 + \frac{\sigma_\epsilon}{\alpha_\epsilon} \hat{\epsilon}\right)$ is
724 evaluated using Eq. 82. In Table 4, we report BPD values for MULAN on CIFAR10 (8M training
725 steps, velocity parameterization) and ImageNet (2M training steps, noise parameterization) using
726 both the VLB-based approach, and the ODE-based approach with $K = 1$ and $K = 20$ importance
727 samples.

728 **I MULAN vs other methods**

729 MULAN is a noise schedule that can be integrated into any diffusion model such as VDM (Kingma
730 et al., 2021), InfoDiffusion (Wang et al., 2023), or i-DODE (Zheng et al., 2023). The shared
731 components among these models are summarized and compared in Table 5.

Table 4: NLL (mean and 95% Confidence Interval for MULAN) on CIFAR10 (8M training steps, velocity parameterization) and ImageNet (2M training steps, noise parameterization) using both the VLB-based approach, and the ODE-based approach. $K = 1$ means that we do not use importance weighted estimator since Eq. 86 is equivalent to Eq. 85 for this case; see Zheng et al. (2023).

Likelihood Estimation type	CIFAR-10 (\downarrow)	Imagenet (\downarrow)
VLB-based	2.59 ± 10^{-3}	3.71 ± 10^{-3}
ODE-based ($K = 1$; Eq. 85)	$2.59 \pm 3 \times 10^{-4}$	3.71 ± 10^{-3}
ODE-based ($K = 20$; Eq. 86)	$2.55 \pm 3 \times 10^{-4}$	3.67 ± 10^{-3}

Table 5: MULAN is a noise schedule that can be integrated into any diffusion model such as VDM (Kingma et al., 2021), InfoDiffusion (Wang et al., 2023), or i-DODE (Zheng et al., 2023). The shared components between MULAN and these models are summarized and compared in this table.

Method	learned noise	multivariate noise	input conditioned noise	auxiliary latents	noise parameterization
VDM (Kingma et al., 2021)	Yes	No	No	No	Monotonic neural network
Blurring Diffusion Model (Hoogeboom & Salimans, 2022)	No	Yes	No	No	Frequency scaling
InfoDiffusion (Wang et al., 2023)	No	No	No	In denoising process	Cosine schedule
MULAN (Ours)	Yes	Yes	Yes	In noising and denoising process	Polynomial, sigmoid