
Elucidating the Design Space of Diffusion-Based Generative Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We argue that the theory and practice of diffusion-based generative models are
2 currently unnecessarily convoluted and seek to remedy the situation by presenting
3 a design space that clearly separates the concrete design choices. This lets us
4 identify several changes to both the sampling and training processes, as well as
5 preconditioning of the score networks. Together, our improvements yield new
6 state-of-the-art FID of 1.79 for CIFAR-10 in a class-conditional setting and 1.97 in
7 an unconditional setting, with much faster sampling (35 network evaluations per
8 image) than prior designs. To further demonstrate their modular nature, we show
9 that our design changes dramatically improve both the efficiency and quality ob-
10 tainable with pre-trained score networks from previous work, including improving
11 the FID of an existing ImageNet-64 model from 2.07 to near-SOTA 1.55.

12 1 Introduction

13 Diffusion-based generative models have emerged as a powerful new framework for neural image
14 synthesis, in both unconditional [15, 32, 41] and conditional [16, 31, 32, 34, 35, 36, 37, 41] settings,
15 even surpassing the quality of GANs [12] in certain situations [9]. They are also rapidly finding use
16 in other domains such as audio [26, 33] and video [18] generation, image segmentation [4, 47] and
17 language translation [30]. As such, there is great interest in applying these models and improving
18 them further in terms of image/distribution quality, training cost, and generation speed.

19 The literature on these models is dense on theory, and derivations of sampling schedule, training
20 dynamics, noise level parameterization, etc., tend to be based as directly as possible on theoretical
21 frameworks, which ensures that the models are on a solid theoretical footing. However, this approach
22 has a danger of obscuring the available design space — a proposed model may appear as a tightly
23 coupled package where no individual component can be modified without breaking the entire system.

24 As our first contribution, we take a look at the theory behind these models from a practical standpoint,
25 focusing more on the “tangible” objects and algorithms that appear in the training and sampling
26 phases, and less on the statistical processes from which they might be derived. The goal is to obtain
27 better insights into how these components are linked together and what degrees of freedom are
28 available in the design of the overall system. We focus on the broad class of models where a neural
29 network is used to model the score [21] of a noise level dependent marginal distribution of the training
30 data corrupted by Gaussian noise. Thus, our work is in the context of *denoising score matching* [44].

31 Our second set of contributions concerns the sampling processes used to synthesize images using
32 diffusion models. We identify the best-performing time discretization for sampling, apply a higher-
33 order Runge–Kutta method for the sampling process, evaluate different sampler schedules, and
34 analyze the usefulness of stochasticity in the sampling process. The result of these improvements is a
35 significant drop in the number of sampling steps required during synthesis, and the improved sampler
36 can be used as a drop-in replacement with several widely used diffusions models [32, 41].

37 The third set of contributions focuses on the training of the score-modeling neural network. While
 38 we continue to rely on the commonly used network architectures (DDPM [15], NCSN [40]), we
 39 provide the first principled analysis of the preconditioning of the networks’ inputs, outputs, and loss
 40 functions in a diffusion model setting and derive best practices for improving the training dynamics.
 41 We also suggest an improved distribution of noise levels during training, and note that non-leaking
 42 augmentation [24] — typically used with GANs — is beneficial for diffusion models as well.

43 Taken together, our contributions enable significant improvements in result quality, e.g., leading
 44 to a record FID of 1.79 on CIFAR-10 [27]. With all key ingredients of the design space explicitly
 45 tabulated, we believe that our approach will allow easier innovation on the individual components,
 46 and thus enable more extensive and targeted exploration of the design space of diffusion models.

47 2 Expressing diffusion models in a common framework

48 Let us denote the data distribution by $p_{\text{data}}(\mathbf{x})$, with standard deviation σ_{data} , and consider the family
 49 of mollified distributions $p(\mathbf{x}; \sigma)$ obtained by adding i.i.d. Gaussian noise of standard deviation σ to
 50 the data. For $\sigma_{\text{max}} \gg \sigma_{\text{data}}$, $p(\mathbf{x}; \sigma_{\text{max}})$ is practically indistinguishable from pure Gaussian noise. The
 51 idea of diffusion models is to randomly sample a noise image $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{max}}^2 \mathbf{I})$, and sequentially
 52 denoise it into images \mathbf{x}_i with noise levels $\sigma_0 = \sigma_{\text{max}} > \sigma_1 > \dots > \sigma_N = 0$ so that at each noise
 53 level $\mathbf{x}_i \sim p(\mathbf{x}_i; \sigma_i)$. The endpoint \mathbf{x}_N of this process is thus distributed according to the data.

54 Song et al. [41] present a stochastic differential equation (SDE) that maintains the desired distribution
 55 p as sample \mathbf{x} evolves over time. This allows the above process to be implemented using a stochastic
 56 solver that both removes and adds noise at each iteration. They also give a corresponding “probability
 57 flow” ordinary differential equation (ODE) where the only source of randomness is the initial noise
 58 image \mathbf{x}_0 . Contrary to the usual order of treatment, we begin by examining the ODE, as it offers a
 59 fruitful setting for analyzing sampling trajectories and their discretizations. The insights carry over to
 60 stochastic sampling, which we reintroduce as a generalization in Section 4.

61 **ODE formulation.** A probability flow ODE [41] continuously increases or reduces noise level of
 62 the image when moving forward or backward in time, respectively. To specify the ODE, we must first
 63 choose a schedule $\sigma(t)$ that defines the desired noise level at time t . For example, setting $\sigma(t) \propto \sqrt{t}$
 64 is mathematically natural, as it corresponds to constant-speed heat diffusion [11]. However, we will
 65 show in Section 3 that the choice of schedule has major practical implications and should not be
 66 made on the basis of theoretical convenience.

67 The defining characteristic of the probability flow ODE is that evolving a sample $\mathbf{x}_a \sim p(\mathbf{x}_a; \sigma(t_a))$
 68 from time t_a to t_b (either forward or backward in time) yields a sample $\mathbf{x}_b \sim p(\mathbf{x}_b; \sigma(t_b))$. Following
 69 previous work [41], this requirement is satisfied (Appendix B) by

$$\mathrm{d}\mathbf{x} = -\dot{\sigma}(t) \sigma(t) \nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma(t)) \mathrm{d}t, \quad (1)$$

70 where the dot denotes a time derivative. $\nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma)$ is the *score function* [21], a vector field that
 71 points towards higher density of data at a given noise level. Intuitively, an infinitesimal forward step
 72 of this ODE nudges the sample away from the data, at a rate that depends on the change in noise level.
 73 Equivalently, a backward step nudges the sample towards the data distribution.

74 **Denosing score matching.** The score function has the remarkable property that it does not depend
 75 on the generally intractable normalization constant of the underlying density function $p(\mathbf{x}; \sigma)$ [21],
 76 and thus can be much easier to evaluate. Specifically, if $D(\mathbf{x}; \sigma)$ is a denoiser function that minimizes
 77 the expected L_2 error for samples drawn from p_{data} separately for every σ , i.e.,

$$\mathbb{E}_{\mathbf{y} \sim p_{\text{data}}} \mathbb{E}_{\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} \|D(\mathbf{y} + \mathbf{n}; \sigma) - \mathbf{y}\|_2^2, \text{ then } \nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma) = (D(\mathbf{x}; \sigma) - \mathbf{x}) / \sigma^2, \quad (2, 3)$$

78 where \mathbf{y} is a training image and \mathbf{n} is noise. In this light, the score function isolates the noise
 79 component from the signal in \mathbf{x} , and Eq. 1 amplifies (or diminishes) it over time. Figure 1 illustrates
 80 the behavior of ideal D in practice. The key observation in diffusion models is that $D(\mathbf{x}; \sigma)$ can be
 81 implemented as a neural network $D_{\theta}(\mathbf{x}; \sigma)$ trained according to Eq. 2. Note that D_{θ} may include
 82 additional pre- and post-processing steps, such as scaling \mathbf{x} to an appropriate dynamic range; we will
 83 return to such *preconditioning* in Section 5.

84 **Time-dependent signal scaling.** Some methods (see Appendix C) introduce an additional scale
 85 schedule $s(t)$ and consider $\mathbf{x} = s(t) \hat{\mathbf{x}}$ to be a scaled version of the original, non-scaled variable

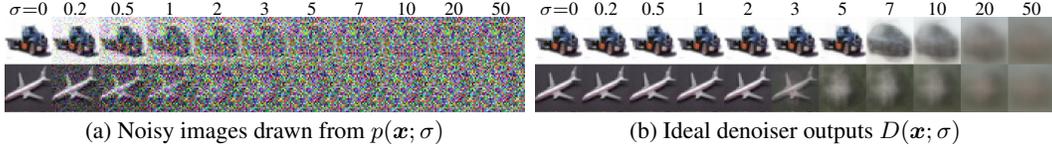


Figure 1: Denoising score matching on CIFAR-10. **(a)** Images from the training set corrupted with varying levels of additive Gaussian noise. High levels of noise lead to oversaturated colors; we normalize the images for cleaner visualization. **(b)** Optimal denoising result from minimizing Eq. 2 analytically (Appendix B). With increasing noise level, the denoised image approaches dataset mean.

Table 1: Specific design choices employed by different model families. N is the number of ODE solver iterations that we wish to execute during sampling. The corresponding sequence of time steps is $\{t_0, t_1, \dots, t_N\}$, where $t_N = 0$. If the model was originally trained for specific choices of N and $\{t_i\}$, the originals are denoted by M and $\{u_j\}$, respectively. The denoiser is defined as $D_\theta(\mathbf{x}; \sigma) = c_{\text{skip}}(\sigma)\mathbf{x} + c_{\text{out}}(\sigma)F_\theta(c_{\text{in}}(\sigma)\mathbf{x}; c_{\text{noise}}(\sigma))$; F_θ represents the raw neural network layers.

	VP [41]	VE [41]	iDDPM [32] + DDIM [39]	Ours
Sampling (Section 3)				
ODE solver	Euler	Euler	Euler	2 nd order Heun
Time steps	$t_{i < N}$	$1 + \frac{i}{N-1}(\epsilon_s - 1)$	$\sigma_{\text{max}}^2 (\sigma_{\text{min}}^2 / \sigma_{\text{max}}^2)^{\frac{i}{N-1}}$	$(\sigma_{\text{max}}^{\frac{1}{\rho}} + \frac{i}{N-1}(\sigma_{\text{min}}^{\frac{1}{\rho}} - \sigma_{\text{max}}^{\frac{1}{\rho}}))^\rho$
Schedule	$\sigma(t)$	$\sqrt{e^{\frac{1}{2}\beta_d t^2 + \beta_{\text{min}} t} - 1}$	\sqrt{t}	t
Scaling	$s(t)$	$1/\sqrt{e^{\frac{1}{2}\beta_d t^2 + \beta_{\text{min}} t}}$	1	1
Network and preconditioning (Section 5)				
Architecture of F_θ	DDPM++	NCSN++	DDPM	(any)
Skip scaling $c_{\text{skip}}(\sigma)$	1	1	1	$\sigma_{\text{data}}^2 / (\sigma^2 + \sigma_{\text{data}}^2)$
Output scaling $c_{\text{out}}(\sigma)$	$-\sigma$	σ	$-\sigma$	$\sigma \cdot \sigma_{\text{data}} / \sqrt{\sigma_{\text{data}}^2 + \sigma^2}$
Input scaling $c_{\text{in}}(\sigma)$	$1/\sqrt{\sigma^2 + 1}$	1	$1/\sqrt{\sigma^2 + 1}$	$1/\sqrt{\sigma^2 + \sigma_{\text{data}}^2}$
Noise cond. $c_{\text{noise}}(\sigma)$	$(M-1)\sigma^{-1}(\sigma)$	$\ln(\frac{1}{2}\sigma)$	$M-1 - \arg \min_j u_j - \sigma $	$\frac{1}{4} \ln(\sigma)$
Training (Section 5)				
Noise distribution	$\sigma^{-1}(\sigma) \sim \mathcal{U}(\epsilon_t, 1)$	$\ln(\sigma) \sim \mathcal{U}(\ln(\sigma_{\text{min}}), \ln(\sigma_{\text{max}}))$	$\sigma = u_j, j \sim \mathcal{U}\{0, M-1\}$	$\ln(\sigma) \sim \mathcal{N}(P_{\text{mean}}, P_{\text{std}}^2)$
Loss weighting $\lambda(\sigma)$	$1/\sigma^2$	$1/\sigma^2$	$1/\sigma^2$ (note: *)	$(\sigma^2 + \sigma_{\text{data}}^2) / (\sigma \cdot \sigma_{\text{data}})^2$
Parameters				
	$\beta_d = 19.9, \beta_{\text{min}} = 0.1$	$\sigma_{\text{min}} = 0.02$	$\bar{\alpha}_j = \sin^2(\frac{\pi}{2} \frac{j}{M(C_2+1)})$	$\sigma_{\text{min}} = 0.002, \sigma_{\text{max}} = 80$
	$\epsilon_s = 10^{-3}, \epsilon_t = 10^{-5}$	$\sigma_{\text{max}} = 100$	$C_1 = 0.001, C_2 = 0.008$	$\sigma_{\text{data}} = 0.5, \rho = 7$
	$M = 1000$		$M = 1000, j_0 = 8^\dagger$	$P_{\text{mean}} = -1.2, P_{\text{std}} = 1.2$

* iDDPM also employs a second loss term L_{vib} † In our tests, $j_0 = 8$ yielded better FID than $j_0 = 0$ used by iDDPM

86 $\hat{\mathbf{x}}$. This changes the time-dependent probability density, and consequently also the ODE solution
87 trajectories. The resulting ODE is a generalization of Eq. 1:

$$d\mathbf{x} = \left[\dot{s}(t) \mathbf{x}/s(t) - s(t)^2 \dot{\sigma}(t) \sigma(t) \nabla_{\mathbf{x}} \log p(\mathbf{x}/s(t); \sigma(t)) \right] dt. \quad (4)$$

88 Note that we explicitly undo the scaling of \mathbf{x} when evaluating the score function to keep the definition
89 of $p(\mathbf{x}; \sigma)$ independent of $s(t)$.

90 **Solution by discretization.** The ODE to be solved is obtained by substituting Eq. 3 into Eq. 4 to
91 define the point-wise gradient, and the solution can be found by numerical integration, i.e., taking
92 finite steps over discrete time intervals. This requires choosing both the integration scheme (e.g.,
93 Euler or a variant of Runge–Kutta), as well as the discrete sampling times $\{t_0, t_1, \dots, t_N\}$. Many
94 prior works rely on Euler’s method, but we show in Section 3 that a 2nd order solver offers a better
95 computational tradeoff. For brevity, we do not provide a separate pseudocode for Euler’s method
96 applied to our ODE here, but it can be extracted from Algorithm 1 by omitting lines 6–8.

97 **Putting it together.** Table 1 presents formulas for reproducing deterministic variants of three
98 earlier methods in our framework. These methods were chosen because they are widely used and
99 achieve state-of-the-art performance, but also because they were derived from different theoretical
100 foundations. Some of our formulas appear quite different from the original papers as indirection and
101 recursion have been removed; see Appendix C for details. The main purpose of this reframing is to

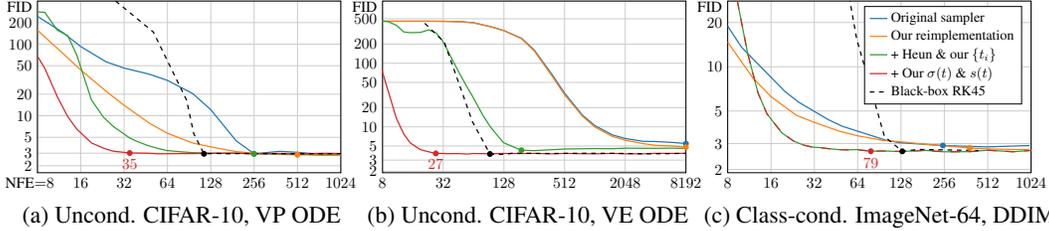


Figure 2: Comparison of deterministic sampling methods using three pre-trained models. For each curve, the dot indicates the lowest NFE whose FID is within 3% of the lowest observed FID.

bring into light all the independent components that often appear tangled together in previous work. In our framework, there are no implicit dependencies between the components — any choices (within reason) for the individual formulas will, in principle, lead to a functioning model.

3 Improvements to deterministic sampling

Our hypothesis is that the choices related to the sampling process are largely independent of the other components, such as network architecture and training details. In other words, the training procedure of D_θ should not dictate $\sigma(t)$, $s(t)$, and $\{t_i\}$, nor vice versa; from the viewpoint of the sampler, D_θ is simply a black box [45, 46]. We test this by evaluating different samplers on three *pre-trained* models, each representing a different theoretical framework and model family. We first measure baseline results for these models using their original sampler implementations, and then bring these samplers into our unified framework using the formulas in Table 1, followed by our improvements.

We evaluate the “DDPM++ cont. (VP)” and “NCSN++ cont. (VE)” models by Song et al. [41] trained on unconditional CIFAR-10 [27] at 32×32 , corresponding to the variance preserving (VP) and variance exploding (VE) formulations [41], originally inspired by DDPM [15] and SMLD [40]. We also evaluate the “ADM (dropout)” model by Dhariwal and Nichol [9] trained on class-conditional ImageNet [8] at 64×64 , corresponding to the improved DDPM (iDDPM) formulation [32]. This model was trained using a discrete set of $M = 1000$ noise levels. Further details are given in Appendix C.

We evaluate the result quality in terms of Fréchet inception distance (FID) [14] computed between 50,000 generated images and all available real images. Figure 2 shows FID as a function of neural function evaluations (NFE), i.e., how many times D_θ is evaluated to produce a single image. Given that the sampling process is dominated entirely by the cost of D_θ , improvements in NFE translate directly to sampling speed. The original deterministic samplers are shown in blue, and the reimplementations of these methods in our unified framework (orange) yield similar but consistently better results. The differences are explained by certain oversights in the original implementations as well as our more careful treatment of discrete noise levels in the case of DDIM; see Appendix C. Note that our reimplementations are fully specified by Algorithm 1 and Table 1, even though the original codebases are structured very differently from each other.

Discretization and higher-order integrators. Solving an ODE numerically is necessarily an approximation of following the true solution trajectory. At each step, the solver introduces *truncation error* that accumulates over the course of N steps. The local error generally scales superlinearly with respect to step size, and thus increasing N improves the accuracy of the solution.

The commonly used Euler’s method is a first order ODE solver with $\mathcal{O}(h^2)$ local error with respect to step size h . Higher-order Runge–Kutta methods [42] scale more favorably but require multiple evaluations of D_θ per step. Through extensive tests, we have found Heun’s 2nd order method [2] (a.k.a. improved Euler, trapezoidal rule) — previously explored in the context of diffusion models by Jolicœur-Martineau et al. [23] — to provide an excellent tradeoff between truncation error and NFE. As illustrated in Algorithm 1, it introduces an additional correction step for \mathbf{x}_{i+1} to account for change in $d\mathbf{x}/dt$ between t_i and t_{i+1} . This correction leads to $\mathcal{O}(h^3)$ local error at the cost of one additional evaluation of D_θ per step. Note that stepping to $\sigma = 0$ would result in a division by zero, so we revert to Euler’s method in this case. We discuss the general family of 2nd order solvers in Appendix D.

The time steps $\{t_i\}$ determine how the step sizes and thus truncation errors are distributed between different noise levels. We provide a detailed analysis in Appendix D, concluding that the step size should decrease monotonically with decreasing σ and it does not need to vary on a per-sample basis.

Algorithm 1 Deterministic sampling using Heun’s 2nd order method with arbitrary $\sigma(t)$ and $s(t)$.

```

1: procedure HEUNSAMPLER( $D_\theta(\mathbf{x}; \sigma)$ ,  $\sigma(t)$ ,  $s(t)$ ,  $t_i \in \{0, \dots, N\}$ )
2:   sample  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \sigma^2(t_0) \mathbf{I})$  ▷ Generate initial sample at  $t_0$ 
3:   for  $i \in \{0, \dots, N - 1\}$  do ▷ Solve Eq. 4 over  $N$  time steps
4:      $\mathbf{d}_i \leftarrow \left( \frac{\dot{\sigma}(t_i)}{\sigma(t_i)} + \frac{\dot{s}(t_i)}{s(t_i)} \right) \mathbf{x}_i - \frac{\dot{\sigma}(t_i) s(t_i)}{\sigma(t_i)} D_\theta \left( \frac{\mathbf{x}_i}{s(t_i)}; \sigma(t_i) \right)$  ▷ Evaluate  $d\mathbf{x}/dt$  at  $t_i$ 
5:      $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + (t_{i+1} - t_i) \mathbf{d}_i$  ▷ Take Euler step from  $t_i$  to  $t_{i+1}$ 
6:     if  $\sigma(t_{i+1}) \neq 0$  then ▷ Apply 2nd order correction unless  $\sigma$  goes to zero
7:        $\mathbf{d}'_i \leftarrow \left( \frac{\dot{\sigma}(t_{i+1})}{\sigma(t_{i+1})} + \frac{\dot{s}(t_{i+1})}{s(t_{i+1})} \right) \mathbf{x}_{i+1} - \frac{\dot{\sigma}(t_{i+1}) s(t_{i+1})}{\sigma(t_{i+1})} D_\theta \left( \frac{\mathbf{x}_{i+1}}{s(t_{i+1})}; \sigma(t_{i+1}) \right)$  ▷ Eval.  $d\mathbf{x}/dt$  at  $t_{i+1}$ 
8:        $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + (t_{i+1} - t_i) \left( \frac{1}{2} \mathbf{d}_i + \frac{1}{2} \mathbf{d}'_i \right)$  ▷ Explicit trapezoidal rule at  $t_{i+1}$ 
9:   return  $\mathbf{x}_N$  ▷ Return noise-free sample at  $t_N$ 

```

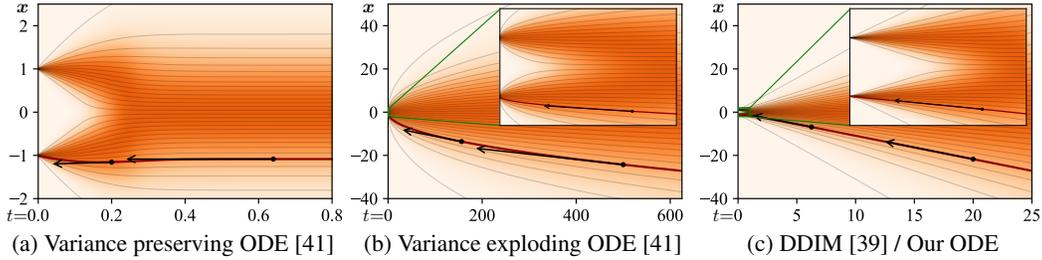


Figure 3: A sketch of ODE curvature in 1D where p_{data} is two Dirac peaks at $\mathbf{x} = \pm 1$. Horizontal t axis is chosen to show $\sigma \in [0, 25]$ in each plot, with insets showing $\sigma \in [0, 1]$ near the data. Example local gradients are shown with black arrows. **(a)** Variance preserving ODE of Song et al. [41] has solution trajectories that flatten out to horizontal lines at large σ . Local gradients start pointing towards data only at small σ . **(b)** Variance exploding variant has extreme curvature near data and the solution trajectories are curved everywhere. **(c)** With the schedule used by DDIM [39] and us, as σ increases the solution trajectories approach straight lines that point towards the mean of data. As $\sigma \rightarrow 0$, the trajectories become linear and point towards the data manifold.

145 We adopt a parameterized scheme where the time steps are defined according to a sequence of noise
146 levels $\{\sigma_i\}$, i.e., $t_i = \sigma^{-1}(\sigma_i)$. We set $\sigma_{i < N} = (Ai + B)^\rho$ and select the constants A and B so that
147 $\sigma_0 = \sigma_{\max}$ and $\sigma_{N-1} = \sigma_{\min}$, which gives

$$\sigma_{i < N} = \left(\sigma_{\max}^{\frac{1}{\rho}} + \frac{i}{N-1} (\sigma_{\min}^{\frac{1}{\rho}} - \sigma_{\max}^{\frac{1}{\rho}}) \right)^\rho \quad \text{and} \quad \sigma_N = 0. \quad (5)$$

148 Here ρ controls how much the steps near σ_{\min} are shortened at the expense of longer steps near σ_{\max} .
149 Our analysis in Appendix D shows that setting $\rho = 3$ nearly equalizes the truncation error at each
150 step, but that ρ in range of 5 to 10 performs much better for sampling images. This suggests that
151 errors near σ_{\min} have a large impact. We set $\rho = 7$ for the remainder of this paper. Results for Heun’s
152 method and Eq. 5 are shown as the green curves in Figure 2. We observe consistent improvement in
153 all cases: Heun’s method reaches the same FID as Euler’s method with considerably lower NFE.

154 **Trajectory curvature and noise schedule.** The shape of the ODE solution trajectories is defined
155 by functions $\sigma(t)$ and $s(t)$. The choice of these functions offers a way to reduce the truncation errors
156 discussed above, as their magnitude can be expected to scale proportional to the curvature of $d\mathbf{x}/dt$.
157 We argue that the best choice for these functions is $\sigma(t) = t$ and $s(t) = 1$, which is also the choice
158 made in DDIM [39]. With this choice, the ODE of Eq. 4 simplifies to $d\mathbf{x}/dt = (\mathbf{x} - D(\mathbf{x}; t))/t$.

159 An immediate consequence is that at any \mathbf{x} and t , a single Euler step to $t = 0$ yields the denoised
160 image $D_\theta(\mathbf{x}; t)$. The tangent of the solution trajectory therefore always points towards the denoiser
161 output. This can be expected to change only slowly with the noise level, which corresponds to largely
162 linear solution trajectories. The 1D ODE sketch of Figure 3c supports this intuition; the solution
163 trajectories approach linear at both large and small noise levels, and have substantial curvature in
164 only a small region in between. The same effect can be seen with real data in Figure 1b, where the
165 change between different denoiser targets occurs in a relatively narrow σ range. With the advocated
166 schedule, this corresponds to high ODE curvature being limited to this same range.

Algorithm 2 Our stochastic sampler with $\sigma(t) = t$ and $s(t) = 1$.

```

1: procedure STOCHASTICSAMPLER( $D_\theta(\mathbf{x}; \sigma)$ ,  $t_i \in \{0, \dots, N\}$ ,  $\gamma_i \in \{0, \dots, N-1\}$ ,  $S_{\text{noise}}$ )
2:   sample  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, t_0^2 \mathbf{I})$ 
3:   for  $i \in \{0, \dots, N-1\}$  do
4:     sample  $\epsilon_i \sim \mathcal{N}(\mathbf{0}, S_{\text{noise}}^2 \mathbf{I})$ 
5:      $\hat{t}_i \leftarrow t_i + \gamma_i t_i$ 
6:      $\hat{\mathbf{x}}_i \leftarrow \mathbf{x}_i + \sqrt{\hat{t}_i^2 - t_i^2} \epsilon_i$ 
7:      $\mathbf{d}_i \leftarrow (\hat{\mathbf{x}}_i - D_\theta(\hat{\mathbf{x}}_i; \hat{t}_i)) / \hat{t}_i$ 
8:      $\mathbf{x}_{i+1} \leftarrow \hat{\mathbf{x}}_i + (t_{i+1} - t_i) \mathbf{d}_i$ 
9:     if  $t_{i+1} \neq 0$  then
10:       $\mathbf{d}'_i \leftarrow (\mathbf{x}_{i+1} - D_\theta(\mathbf{x}_{i+1}; t_{i+1})) / t_{i+1}$ 
11:       $\mathbf{x}_{i+1} \leftarrow \hat{\mathbf{x}}_i + (t_{i+1} - \hat{t}_i) (\frac{1}{2} \mathbf{d}_i + \frac{1}{2} \mathbf{d}'_i)$ 
12:   return  $\mathbf{x}_N$ 

```

$\triangleright \gamma_i = \begin{cases} \min\left(\frac{S_{\text{churn}}}{N}, \sqrt{2}-1\right) & \text{if } t_i \in [S_{\text{min}}, S_{\text{max}}] \\ 0 & \text{otherwise} \end{cases}$
 \triangleright Select temporarily increased noise level \hat{t}_i
 \triangleright Add new noise to move from t_i to \hat{t}_i
 \triangleright Evaluate $d\mathbf{x}/dt$ at \hat{t}_i
 \triangleright Take Euler step from \hat{t}_i to t_{i+1}
 \triangleright Apply 2nd order correction

167 The effect of setting $\sigma(t) = t$ and $s(t) = 1$ is shown as the red curves in Figure 2. As DDIM already
168 employs these same choices, the red curve is identical to the green one for ImageNet-64. However,
169 VP and VE benefit considerably from switching away from their original schedules.

170 **Discussion.** The choices that we made in this section to improve deterministic sampling are
171 summarized in the *Sampling* part of Table 1. Together, they reduce the NFE needed to reach high-
172 quality results by a large factor: $7.3\times$ for VP, $300\times$ for VE, and $3.2\times$ for DDIM, corresponding to
173 the highlighted NFE values in Figure 2. In practice, we can generate 26.3 high-quality CIFAR-10
174 images per second on a single NVIDIA V100. The consistency of improvements corroborates our
175 hypothesis that the sampling process is orthogonal to how each model was originally trained. As
176 further validation, we show results for the adaptive RK45 method [10] using our schedule as the
177 dashed black curves in Figure 2; the cost of this sophisticated ODE solver outweighs its benefits.

178 4 Stochastic sampling

179 Deterministic sampling offers many benefits, e.g., the ability to turn real images into their corre-
180 sponding latent representations by inverting the ODE. However, it tends to lead to worse output
181 quality [39, 41] than stochastic sampling that injects fresh noise into the image in each step. Given
182 that ODEs and SDEs recover the same distributions in theory, what exactly is the role of stochasticity?

183 **Background.** The SDEs of Song et al. [41] can be generalized [19, 48] as a sum of the probability
184 flow ODE of Eq. 1 and a varying-rate *Langevin diffusion* SDE [13]:

$$d\mathbf{x}_\pm = \underbrace{-\dot{\sigma}(t)\sigma(t)\nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma(t)) dt}_{\text{probability flow ODE (Eq. 1)}} \pm \underbrace{\beta(t)\sigma(t)^2 \nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma(t)) dt}_{\text{deterministic noise decay}} + \underbrace{\sqrt{2\beta(t)}\sigma(t) d\omega_t}_{\text{noise injection}}, \quad (6)$$

Langevin diffusion SDE

185 where $d\omega_t$ is the standard Wiener process. $d\mathbf{x}_+$ and $d\mathbf{x}_-$ are now separate SDEs for moving forward
186 and backward in time, related by the time reversal formula of Anderson [1]. The Langevin term can
187 further be seen as a combination of a deterministic score-based denoising term and a stochastic noise
188 injection term, whose net noise level contributions cancel out. As such, $\beta(t)$ effectively expresses the
189 relative rate at which existing noise is replaced with new noise. The SDEs of Song et al. [41] are
190 recovered with the choice $\beta(t) = \dot{\sigma}(t)/\sigma(t)$, whereby the score vanishes from the forward SDE.

191 This perspective reveals why stochasticity is helpful in practice: The implicit Langevin diffusion
192 drives the sample towards the desired marginal distribution at a given time, actively correcting for
193 any errors made in earlier sampling steps. On the other hand, approximating the Langevin term
194 with discrete SDE solver steps introduces error in itself. Previous results [3, 23, 39, 41] suggest that
195 non-zero $\beta(t)$ is helpful, but as far as we can tell, the implicit choice for $\beta(t)$ in Song et al. [41] enjoys
196 no special properties. Hence, the optimal amount of stochasticity should be determined empirically.

197 **Our stochastic sampler.** We propose a stochastic sampler that combines the existing higher-order
198 ODE integrator with explicit Langevin-like “churn” of adding and removing noise. A pseudocode is
199 given in Algorithm 2. At each step i , given the sample \mathbf{x}_i at noise level t_i ($= \sigma(t_i)$), we perform two

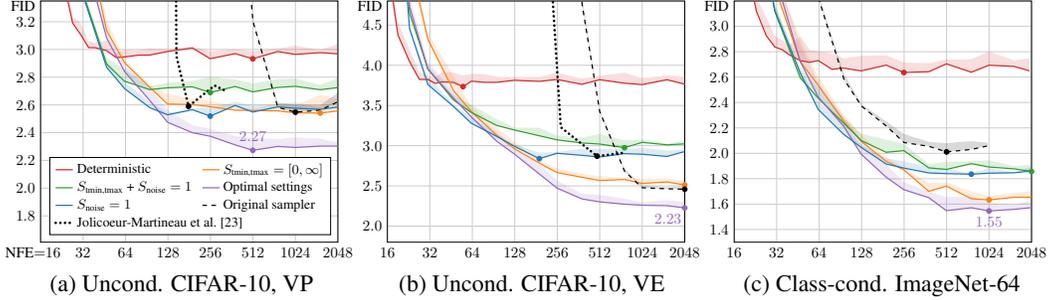


Figure 4: Evaluation of our stochastic sampler (Algorithm 2). The purple curve corresponds to optimal choices for $\{S_{\text{churn}}, S_{\text{tmin}}, S_{\text{tmax}}, S_{\text{noise}}\}$; orange, blue, and green correspond to disabling the effects of $S_{\text{tmin,tmax}}$ and/or S_{noise} . The red curves show reference results for our deterministic sampler (Algorithm 1), equivalent to setting $S_{\text{churn}} = 0$. The dashed black curves correspond to the original stochastic samplers from previous work: Euler–Maruyama [41] for VP, predictor-corrector [41] for VE, and iDDPM [32] for ImageNet-64. The dots indicate lowest observed FID.

200 sub-steps. First, we add noise to the sample according to a factor $\gamma_i \geq 0$ to reach a higher noise level
 201 $\hat{t}_i = t_i + \gamma_i t_i$. Second, from the increased-noise sample \hat{x}_i , we solve the ODE backward from \hat{t}_i to
 202 t_{i+1} with a single step. This yields a sample x_{i+1} with noise level t_{i+1} , and the iteration continues.
 203 We stress that this is not a general-purpose SDE solver, but a sampling procedure tailored for the
 204 specific problem. Its correctness stems from the alternation of two sub-steps that each maintain the
 205 correct distribution (up to truncation error in the ODE step).

206 **Practical considerations.** We have observed (see Appendix E) that excessive Langevin-like addition
 207 and removal of noise results in gradual loss of detail in the generated images with all datasets and
 208 denoiser networks. There is also a drift toward oversaturated colors at very low and high noise levels.
 209 We suspect that practical denoisers induce a slightly non-conservative vector field in Eq. 3, violating
 210 the premises of Langevin diffusion and causing these detrimental effects. Notably, our experiments
 211 with analytical denoisers (such as the one in Figure 1b) have not shown such degradation.

212 If the degradation is caused by flaws in $D_\theta(x; \sigma)$, they can only be remedied using heuristic means
 213 during sampling. We address the drift toward oversaturated colors by only enabling stochasticity
 214 within a specific range of noise levels $t_i \in [S_{\text{tmin}}, S_{\text{tmax}}]$. For these noise levels, we define $\gamma_i =$
 215 S_{churn}/N , where S_{churn} controls the overall amount of stochasticity. We further clamp γ_i to never
 216 introduce more new noise than what is already present in the image. Finally, we have found that the
 217 loss of detail can be partially counteracted by setting S_{noise} slightly above 1 to inflate the standard
 218 deviation for the newly added noise. This suggests that a major component of the hypothesized
 219 non-conservativity of $D_\theta(x; \sigma)$ is a tendency to remove slightly too much noise — most likely due to
 220 regression toward the mean that can be expected to happen with any L_2 -trained denoiser [28].

221 **Evaluation.** Figure 4 shows that our stochastic sampler outperforms previous samplers [32, 41]
 222 by a significant margin, especially at low step counts. In particular, through sampler improvements
 223 alone, we are able to bring the ImageNet-64 model that originally achieved FID 2.07 [9] to 1.55
 224 that is very close to the state-of-the-art; previously, FID 1.48 has been reported for cascaded diffusion
 225 [16], 1.55 for classifier-free guidance [17], and 1.52 for StyleGAN-XL [38]. While our results
 226 showcase the potential gains achievable through sampler improvements, they also highlight the main
 227 shortcoming of stochasticity: For best results, one must make several heuristic choices — either
 228 implicit or explicit — that depend on the specific model. Indeed, we had to find the optimal values of
 229 $\{S_{\text{churn}}, S_{\text{tmin}}, S_{\text{tmax}}, S_{\text{noise}}\}$ on a case-by-case basis using grid search (Appendix E). This raises a general
 230 concern that using stochastic sampling as the primary means of evaluating model improvements
 231 may inadvertently end up influencing the design choices related to model architecture and training.

232 5 Preconditioning and training

233 There are various known good practices for training neural networks in a supervised fashion. For
 234 example, it is advisable to keep input and output signal magnitudes fixed to, e.g., unit variance, and to
 235 avoid large variation in gradient magnitudes on a per-sample basis [5, 20]. Training a neural network
 236 to model D directly would be far from ideal — for example, as the input $x = y + n$ is a combination

Table 2: Evaluation of our training improvements. The starting point (config A) is VP & VE using our **deterministic** sampler. At the end (configs E,F), VP & VE only differ in the architecture of F_θ .

Training configuration	CIFAR-10 [27] at 32×32				FFHQ [25] 64×64		AFHQv2 [7] 64×64	
	Conditional		Unconditional		Unconditional		Unconditional	
	VP	VE	VP	VE	VP	VE	VP	VE
A Baseline [41] (*pre-trained)	2.48	3.11	3.01*	3.77*	3.39	25.95	2.58	18.52
B + Adjust hyperparameters	2.18	2.48	2.51	2.94	3.13	22.53	2.43	23.12
C + Redistribute capacity	2.08	2.52	2.31	2.83	2.78	41.62	2.54	15.04
D + Our preconditioning	2.09	2.64	2.29	3.10	2.94	3.39	2.79	3.81
E + Our loss function	1.88	1.86	2.05	1.99	2.60	2.81	2.29	2.28
F + Non-leaky augmentation	1.79	1.79	1.97	1.98	2.39	2.53	1.96	2.16
NFE	35	35	35	35	79	79	79	79

of clean signal \mathbf{y} and noise $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, its magnitude varies immensely depending on noise level σ . For this reason, the common practice is to not represent D_θ as a neural network directly, but instead train a different network F_θ from which D_θ is derived.

Previous methods [32, 39, 41] address the input scaling via a σ -dependent normalization factor and attempt to precondition the output by training F_θ to predict \mathbf{n} scaled to unit variance, from which the signal is then reconstructed via $D_\theta(\mathbf{x}; \sigma) = \mathbf{x} - \sigma F_\theta(\cdot)$. This has the drawback that at large σ , the network needs to fine-tune its output carefully to cancel out the existing noise \mathbf{n} exactly and give the output at the correct scale; note that any errors made by the network are amplified by a factor of σ . In this situation, it would seem much easier to predict the expected output $D(\mathbf{x}; \sigma)$ directly. To this end, we propose to precondition the neural network with a σ -dependent skip connection that allows it to estimate either \mathbf{y} or \mathbf{n} , or something in between. We thus write D_θ in the following form:

$$D_\theta(\mathbf{x}; \sigma) = c_{\text{skip}}(\sigma) \mathbf{x} + c_{\text{out}}(\sigma) F_\theta(c_{\text{in}}(\sigma) \mathbf{x}; c_{\text{noise}}(\sigma)), \quad (7)$$

where F_θ is the neural network to be trained, $c_{\text{skip}}(\sigma)$ modulates the skip connection, $c_{\text{in}}(\sigma)$ and $c_{\text{out}}(\sigma)$ scale the input and output magnitudes, and $c_{\text{noise}}(\sigma)$ maps noise level σ into a conditioning input for F_θ . Taking a weighted expectation of Eq. 2 over the noise levels gives the overall training loss $\mathbb{E}_{\sigma, \mathbf{y}, \mathbf{n}} [\lambda(\sigma) \|D(\mathbf{y} + \mathbf{n}; \sigma) - \mathbf{y}\|_2^2]$, where $\sigma \sim p_{\text{train}}$, $\mathbf{y} \sim p_{\text{data}}$, and $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. The probability of sampling a given noise level σ is given by $p_{\text{train}}(\sigma)$ and the corresponding weight is given by $\lambda(\sigma)$. We can equivalently express this loss with respect to the raw network output F_θ in Eq. 7:

$$\mathbb{E}_{\sigma, \mathbf{y}, \mathbf{n}} \left[\underbrace{\lambda(\sigma) c_{\text{out}}(\sigma)^2}_{\text{effective weight}} \left\| \underbrace{F_\theta(c_{\text{in}}(\sigma) \cdot (\mathbf{y} + \mathbf{n}); c_{\text{noise}}(\sigma))}_{\text{network output}} - \underbrace{\frac{1}{c_{\text{out}}(\sigma)} (\mathbf{y} - c_{\text{skip}}(\sigma) \cdot (\mathbf{y} + \mathbf{n}))}_{\text{effective training target}} \right\|_2^2 \right]. \quad (8)$$

This form reveals the effective training target of F_θ , allowing us to determine suitable choices for the preconditioning functions from first principles. As detailed in Appendix B, we derive our choices shown in Table 1 by requiring network inputs and training targets to have unit variance (c_{in} , c_{out}), and amplifying errors in F_θ as little as possible (c_{skip}). The formula for c_{noise} is chosen empirically.

Table 2 shows FID for a series of training setups, evaluated using our deterministic sampler from Section 3. We start with the baseline training setup of Song et al. [41], which differs considerably between the VP and VE cases; we provide separate results for each (config A). To obtain a more meaningful point of comparison, we re-adjust the basic hyperparameters (config B) and improve the expressive power of the model (config C) by removing the lowest-resolution layers and doubling the capacity of the highest-resolution layers instead; see Appendix F for further details. We then replace the original choices of $\{c_{\text{in}}, c_{\text{out}}, c_{\text{noise}}, c_{\text{skip}}\}$ with our preconditioning (config D), which keeps the results largely unchanged — except for VE that improves considerably at 64×64 resolution. Instead of improving FID per se, the main benefit of our preconditioning is that it makes the training more robust, enabling us to turn our focus on redesigning the loss function without adverse effects.

Loss weighting and sampling. Eq. 8 shows that training F_θ as preconditioned in Eq. 7 incurs an effective per-sample loss weight of $\lambda(\sigma) c_{\text{out}}(\sigma)^2$. To balance the effective loss weights, we set $\lambda(\sigma) = 1/c_{\text{out}}(\sigma)^2$, which also equalizes the initial training loss over the entire σ range as shown in Figure 5a (green curve). Finally, we need to select $p_{\text{train}}(\sigma)$, i.e., how to choose noise levels during training. Inspecting the per- σ loss after training (blue and orange curves) reveals that a significant reduction is possible only at intermediate noise levels; at very low levels, it is both difficult and

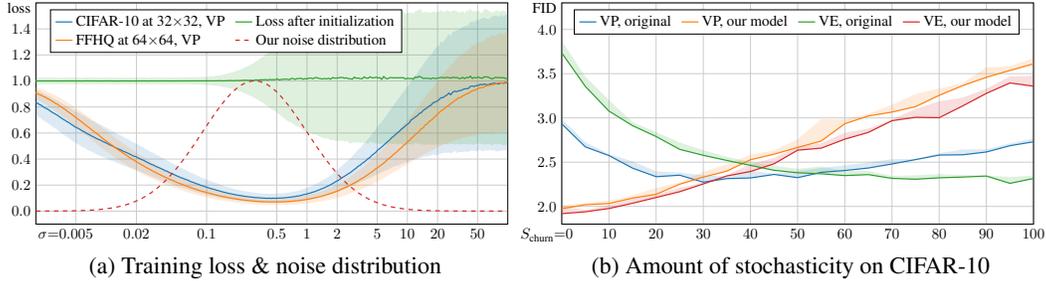


Figure 5: **(a)** Observed initial (green) and final loss per noise level, representative of the the 32×32 (blue) and 64×64 (orange) models considered in this paper. The shaded regions represent the standard deviation over 10k random samples. Our proposed training sample density is shown by the dashed red curve. **(b)** For the original training setup of Song et al. [41], stochastic sampling is highly beneficial (blue, green), while deterministic sampling ($S_{\text{churn}} = 0$) leads to relatively poor FID. For our training setup, the situation is reversed (orange, red); stochastic sampling is not only unnecessary but harmful.

274 irrelevant to discern the vanishingly small noise component, whereas at high levels the training targets
 275 are always dissimilar from the correct answer that approaches dataset average. Therefore, we target
 276 the training efforts to the relevant range using a simple log-normal distribution for $p_{\text{train}}(\sigma)$ as detailed
 277 in Table 1 and illustrated in Figure 5a (red curve).

278 Table 2 shows that our proposed p_{train} and λ (config E) lead to a dramatic improvement in FID
 279 in all cases when used in conjunction with our preconditioning (config D). In concurrent work,
 280 Choi et al. [6] propose a similar scheme to prioritize noise levels that are most relevant w.r.t. forming
 281 the perceptually recognizable content of the image. However, they only consider the choice of λ in
 282 isolation, which results in a smaller overall improvement.

283 **Augmentation regularization.** To prevent potential overfitting that often plagues diffusion models
 284 with smaller datasets, we borrow an augmentation pipeline from the GAN literature [24]. The pipeline
 285 consists of various geometric transformations (see Appendix F) that we apply to a training image
 286 prior to adding noise. To prevent the augmentations from leaking to the generated images, we provide
 287 the augmentation parameters as a conditioning input to F_{θ} ; during inference we set the them to zero
 288 to guarantee that only non-augmented images are generated. Table 2 shows that data augmentation
 289 provides a consistent improvement (config F) that yields new state-of-the-art FIDs of 1.79 and 1.97
 290 for conditional and unconditional CIFAR-10, beating the previous records of 1.85 [38] and 2.10 [43].

291 6 Conclusions

292 Our approach of putting diffusion models to a common framework exposes a modular design. This
 293 allows a targeted investigation of individual components, potentially helping to better cover the
 294 viable design space. In our tests this let us to simply replace the samplers in various earlier models,
 295 drastically improving the results. For example, in ImageNet-64 our sampler turned an average
 296 model (FID 2.07) to a challenger (1.55) for the current SOTA model (1.48) [16]. We also obtained
 297 new state-of-the-art results on CIFAR-10 while using only 35 model evaluations, deterministic
 298 sampling, and a small network. The current high-resolution diffusion models rely either on separate
 299 super-resolution steps [16, 31, 35], subspace projection [22], very large networks [9, 41], or hybrid
 300 approaches [34, 36, 43] — we believe that our contributions are orthogonal to these extensions. That
 301 said, many of our parameter values may need to be re-adjusted for higher resolution datasets.

302 Interestingly, the relevance of stochastic sampling appears to diminish as the model itself improves,
 303 as shown in Figure 5b. Intuitively the role of stochasticity is to correct approximation errors, and
 304 if the approximation errors are minimal to begin with, there is nothing left to do. Nevertheless, we
 305 feel that the precise interaction between stochasticity and the training objective remains a promising
 306 avenue for future work, and stochasticity likely continues to be beneficial for more diverse datasets.

307 **Negative societal impacts** Our advances in sample quality can potentially amplify negative societal
 308 effects when used in a large-scale system like DALL·E 2, including types of disinformation or
 309 emphasizing stereotypes and harmful biases [29]. The training and sampling of diffusion models needs
 310 a lot of electricity; our project consumed $\sim 250\text{MWh}$ on an in-house cluster of NVIDIA V100s.

311 **References**

- 312 [1] B. D. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*,
313 12(3):313–326, 1982.
- 314 [2] U. M. Ascher and L. R. Petzold. *Computer Methods for Ordinary Differential Equations and Differential-*
315 *Algebraic Equations*. Society for Industrial and Applied Mathematics, 1998.
- 316 [3] F. Bao, C. Li, J. Zhu, and B. Zhang. Analytic-DPM: an analytic estimate of the optimal reverse variance in
317 diffusion probabilistic models. In *Proc. ICLR*, 2022.
- 318 [4] D. Baranchuk, A. Voynov, I. Rubachev, V. Khruikov, and A. Babenko. Label-efficient semantic segmenta-
319 tion with diffusion models. In *Proc. ICLR*, 2022.
- 320 [5] C. M. Bishop. *Neural networks for pattern recognition*. Oxford University Press, USA, 1995.
- 321 [6] J. Choi, J. Lee, C. Shin, S. Kim, H. Kim, and S. Yoon. Perception prioritized training of diffusion models.
322 *CoRR*, abs/2204.00227, 2022.
- 323 [7] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha. StarGAN v2: Diverse image synthesis for multiple domains. In *Proc.*
324 *CVPR*, 2020.
- 325 [8] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image
326 database. In *Proc. CVPR*, pages 248–255. IEEE, 2009.
- 327 [9] P. Dhariwal and A. Q. Nichol. Diffusion models beat GANs on image synthesis. In *Proc. NeurIPS*, 2021.
- 328 [10] J. R. Dormand and P. J. Prince. A family of embedded Runge-Kutta formulae. *Journal of computational*
329 *and applied mathematics*, 6(1):19–26, 1980.
- 330 [11] J. B. J. Fourier, G. Darboux, et al. *Théorie analytique de la chaleur*, volume 504. Didot Paris, 1822.
- 331 [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio.
332 Generative adversarial networks. In *Proc. NIPS*, 2014.
- 333 [13] U. Grenander and M. I. Miller. Representations of knowledge in complex systems. *Journal of the Royal*
334 *Statistical Society: Series B (Methodological)*, 56(4):549–581, 1994.
- 335 [14] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. GANs trained by a two time-scale
336 update rule converge to a local Nash equilibrium. In *Proc. NIPS*, 2017.
- 337 [15] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *Proc. NeurIPS*, 2020.
- 338 [16] J. Ho, C. Saharia, W. Chan, D. J. Fleet, M. Norouzi, and T. Salimans. Cascaded diffusion models for high
339 fidelity image generation. *Journal of Machine Learning Research*, 23, 2022.
- 340 [17] J. Ho and T. Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative*
341 *Models and Downstream Applications*, 2021.
- 342 [18] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet. Video diffusion models. *CoRR*,
343 abs/2204.03458, 2022.
- 344 [19] C.-W. Huang, J. H. Lim, and A. C. Courville. A variational perspective on diffusion-based generative
345 models and score matching. In *Proc. NeurIPS*, 2021.
- 346 [20] L. Huang, J. Qin, Y. Zhou, F. Zhu, L. Liu, and L. Shao. Normalization techniques in training DNNs:
347 Methodology, analysis and application. *CoRR*, abs/2009.12836, 2020.
- 348 [21] A. Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine*
349 *Learning Research*, 6(24):695–709, 2005.
- 350 [22] B. Jing, G. Corso, R. Berlinghieri, and T. Jaakkola. Subspace diffusion generative models. *CoRR*,
351 abs/2205.01490, 2022.
- 352 [23] A. Jolicœur-Martineau, K. Li, R. Piché-Taillefer, T. Kachman, and I. Mitliagkas. Gotta go fast when
353 generating data with score-based models. *CoRR*, abs/2105.14080, 2021.
- 354 [24] T. Karras, M. Aittala, J. Hellsten, S. Laine, J. Lehtinen, and T. Aila. Training generative adversarial
355 networks with limited data. In *Proc. NeurIPS*, 2020.
- 356 [25] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks.
357 In *Proc. CVPR*, 2018.
- 358 [26] Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro. DiffWave: A versatile diffusion model for audio
359 synthesis. In *Proc. ICLR*, 2021.
- 360 [27] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of
361 Toronto, 2009.
- 362 [28] J. Lehtinen, J. Munkberg, J. Hasselgren, S. Laine, T. Karras, M. Aittala, and T. Aila. Noise2Noise:
363 Learning image restoration without clean data. In *Proc. ICML*, 2018.
- 364 [29] P. Mishkin, L. Ahmad, M. Brundage, G. Krueger, and G. Sastry. DALL·E 2 preview – risks and limitations.
365 *OpenAI*, 2022.
- 366 [30] E. Nachmani and S. Dovrat. Zero-shot translation using diffusion models. *CoRR*, abs/2111.01471, 2021.
- 367 [31] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen.
368 GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. *CoRR*,
369 abs/2112.10741, 2021.
- 370 [32] A. Q. Nichol and P. Dhariwal. Improved denoising diffusion probabilistic models. In *Proc. ICML*, volume
371 139, pages 8162–8171, 2021.

- 372 [33] V. Popov, I. Vovk, V. Gogoryan, T. Sadekova, and M. Kudinov. Grad-TTS: A diffusion probabilistic model
373 for text-to-speech. In *Proc. ICML*, volume 139, pages 8599–8608, 2021.
- 374 [34] K. Preechakul, N. Chatthee, S. Wizadwongsa, and S. Suwajanakorn. Diffusion autoencoders: Toward a
375 meaningful and decodable representation. In *Proc. CVPR*, 2022.
- 376 [35] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. Hierarchical text-conditional image generation
377 with CLIP latents. Technical report, OpenAI, 2022.
- 378 [36] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with
379 latent diffusion models. In *Proc. CVPR*, 2022.
- 380 [37] C. Saharia, W. Chan, H. Chang, C. A. Lee, J. Ho, T. Salimans, D. J. Fleet, and M. Norouzi. Palette:
381 Image-to-image diffusion models. *CoRR*, abs/2111.05826, 2021.
- 382 [38] A. Sauer, K. Schwarz, and A. Geiger. StyleGAN-XL: Scaling StyleGAN to large diverse datasets. *CoRR*,
383 abs/2201.00273, 2022.
- 384 [39] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. In *Proc. ICLR*, 2021.
- 385 [40] Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution. In *Proc.*
386 *NeurIPS*, 2019.
- 387 [41] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative
388 modeling through stochastic differential equations. In *Proc. ICLR*, 2021.
- 389 [42] E. Süli and D. Mayers. *An Introduction to Numerical Analysis*. Cambridge University Press, 2003.
- 390 [43] A. Vahdat, K. Kreis, and J. Kautz. Score-based generative modeling in latent space. In *Proc. NeurIPS*,
391 2021.
- 392 [44] P. Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*,
393 23(7):1661–1674, 2011.
- 394 [45] D. Watson, W. Chan, J. Ho, and M. Norouzi. Learning fast samplers for diffusion models by differentiating
395 through sample quality. In *Proc. ICLR*, 2022.
- 396 [46] D. Watson, J. Ho, M. Norouzi, and W. Chan. Learning to efficiently sample from diffusion probabilistic
397 models. *CoRR*, abs/2106.03802, 2021.
- 398 [47] J. Wolleb, R. Sandkühler, F. Bieder, P. Valmaggia, and P. C. Cattin. Diffusion models for implicit image
399 segmentation ensembles. In *Medical Imaging with Deep Learning*, 2022.
- 400 [48] Q. Zhang and Y. Chen. Diffusion normalizing flow. In *Proc. NeurIPS*, 2021.

401 **Checklist**

- 402 1. For all authors...
- 403 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
404 contributions and scope? [Yes]
- 405 (b) Did you describe the limitations of your work? [Yes] Section 6. The main limitations
406 of the analysis relate to the set of tested datasets and their limited resolution.
- 407 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Section 6.
- 408 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
409 them? [Yes]
- 410 2. If you are including theoretical results...
- 411 (a) Did you state the full set of assumptions of all theoretical results? [No] We follow
412 common application-specific assumptions about the probability distributions, functions
413 and other components, but do not exhaustively specify them, or consider pathological
414 corner cases.
- 415 (b) Did you include complete proofs of all theoretical results? [No] Our equations and
416 algorithms build on previously known results, and highlight their practical aspects
417 through mostly readily verifiable algebraic manipulations (Appendix B). We do not
418 explicitly present all details of the derivations, and assume that the previous results are
419 sufficiently rigorously proven in the respective literature.
- 420 3. If you ran experiments...
- 421 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
422 perimental results (either in the supplemental material or as a URL)? [Yes] Code is
423 provided in the supplemental material. We will also release it in public, along with
424 pre-trained models.
- 425 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
426 were chosen)? [Yes] Appendix F.
- 427 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
428 ments multiple times)? [Yes] Shaded regions in Figures 4 and 5.
- 429 (d) Did you include the total amount of compute and the type of resources used (e.g., type
430 of GPUs, internal cluster, or cloud provider)? [Yes] Section 6.
- 431 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 432 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 433 (b) Did you mention the license of the assets? [Yes] Appendix F.
- 434 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 435 (d) Did you discuss whether and how consent was obtained from people whose data you’re
436 using/curating? [N/A]
- 437 (e) Did you discuss whether the data you are using/curating contains personally identifiable
438 information or offensive content? [N/A]
- 439 5. If you used crowdsourcing or conducted research with human subjects...
- 440 (a) Did you include the full text of instructions given to participants and screenshots, if
441 applicable? [N/A]
- 442 (b) Did you describe any potential participant risks, with links to Institutional Review
443 Board (IRB) approvals, if applicable? [N/A]
- 444 (c) Did you include the estimated hourly wage paid to participants and the total amount
445 spent on participant compensation? [N/A]