

LEARNING FAST SAMPLERS FOR DIFFUSION MODELS BY DIFFERENTIATING THROUGH SAMPLE QUALITY

Daniel Watson*, William Chan, Jonathan Ho & Mohammad Norouzi

Google Research, Brain Team

{watsondaniel, williamchan, jonathanho, mnorouzi}@google.com

ABSTRACT

Diffusion models have emerged as an expressive family of generative models rivaling GANs in sample quality and autoregressive models in likelihood scores. Standard diffusion models typically require hundreds of forward passes through the model to generate a single high-fidelity sample. We introduce Differentiable Diffusion Sampler Search (DDSS): a method that optimizes fast samplers for any pre-trained diffusion model by differentiating through sample quality scores. We present Generalized Gaussian Diffusion Models (GGDM), a family of flexible non-Markovian samplers for diffusion models. We show that optimizing the degrees of freedom of GGDM samplers by maximizing sample quality scores via gradient descent leads to improved sample quality. Our optimization procedure backpropagates through the sampling process using the reparametrization trick and gradient rematerialization. DDSS achieves strong results on unconditional image generation across various datasets (*e.g.*, FID scores on LSUN church 128x128 of 11.6 with only 10 inference steps, and 4.82 with 20 steps, compared to 51.1 and 14.9 with strongest DDPM/DDIM baselines). Our method is compatible with any pre-trained diffusion model without fine-tuning or re-training required.

1 INTRODUCTION

Denosing Diffusion Probabilistic Models (DDPM) (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020) have emerged as a powerful family of generative models, capable of synthesizing high-quality images, audio, and 3D shapes (Ho et al., 2020; 2021; Chen et al., 2021a;b; Cai et al., 2020; Luo & Hu, 2021). Recent work (Dhariwal & Nichol, 2021; Ho et al., 2021) shows that DDPMs can outperform Generative Adversarial Networks (GAN) (Goodfellow et al., 2014; Brock et al., 2018) in generation quality, but unlike GANs, DDPMs admit likelihood computation and much more stable training dynamics (Arjovsky et al., 2017; Gulrajani et al., 2017).

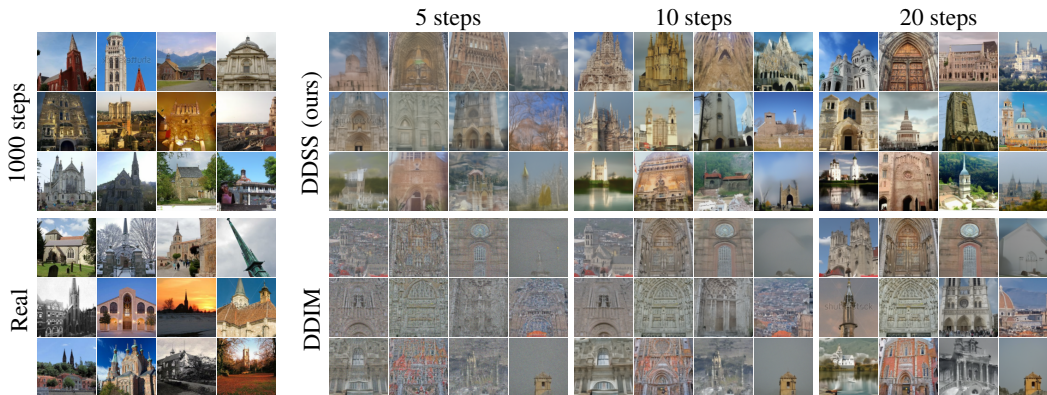


Figure 1: Non-cherry-picked samples from DDSS (ours) and strongest DDIM($\eta = 0$) baseline for unconditional DDPMs trained on LSUN churches 128x128. All samples are generated with the same random seed. Original DDPM samples (1000 steps) and training images are shown on the left.

*Work done as part of the Google AI Residency.

However, GANs are typically much more efficient than DDPMs at generation time, often requiring a single forward pass through the generator network, whereas DDPMs require hundreds of forward passes through a U-Net model. Instead of learning a generator directly, DDPMs learn to convert noisy data to less noisy data starting from pure noise, which leads to a wide variety of feasible strategies for sampling (Song et al., 2021b). In particular, at inference time, DDPMs allow controlling the number of forward passes (a.k.a. *inference steps*) through the denoising network (Song et al., 2020; Nichol & Dhariwal, 2021).

It has been shown both empirically and mathematically that, for any sufficiently good DDPM, more inference steps leads to better log-likelihood and sample quality (Nichol & Dhariwal, 2021; Kingma et al., 2021). In practice, the minimum number of inference steps to achieve competitive sample quality is highly problem-dependent, *e.g.*, depends on the complexity of the dataset, and the strength of the conditioning signal if the task is conditional. Given the importance of generation speed, recent work (Song et al., 2020; Chen et al., 2021a; Watson et al., 2021) has explored reducing the number of steps required for high quality sampling with pretrained diffusion models. See Section 7 for a more complete review of prior work on few-step sampling.

This paper treats the design of fast samplers for diffusion models as a differentiable optimization problem, and proposes *Differentiable Diffusion Sampler Search* (DDSS). Our key observation is that one can unroll the sampling chain of a diffusion model and use reparametrization trick (Kingma & Welling, 2013) and gradient rematerialization (Kumar et al., 2019a) to optimize over a class of parametric few-step samplers with respect to a global objective function. Our class of parametric samplers, which we call Generalized Gaussian Diffusion Model (GGDM), includes Denoising Diffusion Implicit Models (DDIM) (Song et al., 2020) as a special case and is motivated by the success of DDIM on fast sampling of diffusion models.

An important challenge for fast DDPM sampling is the mismatch between the training objective (*e.g.*, ELBO or weighted ELBO) and sample quality. Prior work (Watson et al., 2021; Song et al., 2021a) finds that samplers that are optimal with respect to ELBO often lead to worse sample quality and Fréchet Inception Distance (FID) scores (Heusel et al., 2017), especially with few inference steps. We propose the use of a *perceptual* loss within the DDSS framework to find high-fidelity diffusion samplers, motivated by prior work showing that their optimization leads to solutions that correlate better with human perception of quality. We empirically find that using DDSS with the Kernel Inception Distance (KID) (Bińkowski et al., 2018) as the perceptual loss indeed leads to fast samplers with significantly better image quality than prior work (see Figure 1). Moreover, our method is robust to different choices of kernels for KID.

Our main contributions are as follows:

1. We propose Differentiable Diffusion Sampler Search (DDSS), which uses the reparametrization trick and gradient rematerialization to optimize over a parametric family of fast samplers for diffusion models.
2. We identify a parametric family of Generalized Gaussian Diffusion Model (GGDM) that admits high-fidelity fast samplers for diffusion models.
3. We show that using DDSS to optimize samplers by minimizing the Kernel Inception Distance leads to fast diffusion model samplers with state-of-the-art sample quality scores.

2 BACKGROUND ON DENOISING DIFFUSION IMPLICIT MODELS

We start with a brief review on DDPM (Ho et al., 2020) and DDIM (Song et al., 2020). DDPMs pre-specify a Markovian forward diffusion process, which gradually adds noise to data in T steps. Following the notation of Ho et al. (2020),

$$q(\mathbf{x}_0, \dots, \mathbf{x}_T) = q(\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (1)$$

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t | \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}), \quad (2)$$

where $q(\mathbf{x}_0)$ is the data distribution, and β_t is the variance of Gaussian noise added at step t . To be able to gradually convert noise to data, DDPMs learn to invert (1) with a model $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$, which

is trained by maximizing a (possibly reweighted) evidence lower bound (ELBO):

$$\mathbb{E}_q \left[D_{\text{KL}}[q(\mathbf{x}_T|\mathbf{x}_0)||p(\mathbf{x}_T)] + \sum_{t=2}^T D_{\text{KL}}[q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)] - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \right] \quad (3)$$

DDPMs specifically choose the model to be parametrized as

$$\begin{aligned} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) &= q \left(\mathbf{x}_{t-1} \left| \mathbf{x}_t, \hat{\mathbf{x}}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)) \right. \right) \\ &= \mathcal{N} \left(\mathbf{x}_{t-1} \left| \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right), \frac{1 - \bar{\alpha}_t}{1 - \bar{\alpha}_{t-1}} \beta_t \mathbf{I}_d \right. \right) \end{aligned} \quad (4)$$

where $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$ for each t . With this parametrization, maximizing the ELBO is equivalent to minimizing a weighted sum of denoising score matching objectives (Vincent, 2011).

The seminal work of Song et al. (2020) presents Denoising Diffusion Implicit Models (DDIM): a family of evidence lower bounds (ELBOs) with corresponding forward diffusion processes and samplers. All of these ELBOs share the same marginals as DDPM, but allow arbitrary choices of posterior variances. Specifically, Song et al. (2020) note that it is possible to construct alternative ELBOs with only a subsequence of the original timesteps $S \subset \{1, \dots, T\}$ that shares the same marginals as the construction above (i.e., $q_S(\mathbf{x}_t|\mathbf{x}_0) = q(\mathbf{x}_t|\mathbf{x}_0)$ for every $t \in S$, so q_S defines a faster sampler compatible with the pre-trained model) by simply using the new contiguous timesteps in the equations above. They also show it is also possible to construct an *infinite* family of non-Markovian processes $\{q_\sigma : \sigma \in [0, 1]^{T-1}\}$ where each q_σ also shares the same marginals as the original forward process with:

$$q_\sigma(\mathbf{x}_0, \dots, \mathbf{x}_T) = q(\mathbf{x}_0)q(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=1}^{T-1} q_\sigma(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{x}_0) \quad (5)$$

and where the posteriors are defined as

$$q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_{t-1} \left| \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 \mathbf{I}_d \right. \right) \quad (6)$$

Song et al. (2020) empirically find that the extreme case of using all-zero variances (*a.k.a.* DDIM($\eta = 0$)) consistently helps with sample quality in the few-step regime. Combined with a good selection of timesteps to evaluate the modeled score function (*a.k.a.* *strides*), DDIM($\eta = 0$) establishes the current state-of-the-art for few-step diffusion model sampling with the smallest inference step budgets. Our key contribution that allows improving sample quality significantly by optimizing sampler families is constructing a family that generalizes DDIM. See Section 4 for a more complete treatment of our novel GGDM family.

3 DIFFERENTIABLE DIFFUSION SAMPLER SEARCH (DDSS)

We now describe DDSS, our approach to search for fast high-fidelity samplers with a limited budget of $K < T$ steps. Our key observation is that one can backpropagate through the sampling process of a diffusion model via the reparamterization trick (Kingma & Welling, 2013). Equipped with this, we can now use stochastic gradient descent to learn fast samplers by optimizing any given differentiable loss function over a minibatch of model samples.

We begin with a pre-trained DDPM and a family of K -step samplers that we wish to optimize for the given DDPM. We parametrize this family’s degrees of freedom as simple transformations of trainable variables. We experiment with the following families in this paper, but emphasize that DDSS is applicable to any other family where model samples are differentiable with respect to the trainable variables:

- **DDIM**: we parametrize the posterior variances with $\sigma_t = \sqrt{1 - \bar{\alpha}_{t-1}} \text{sigmoid}(v_t)$, where v_1, \dots, v_K are trainable variables (the $\sqrt{1 - \bar{\alpha}_{t-1}}$ constant ensures real-valued mean coefficients; see the square root in Equation 6).

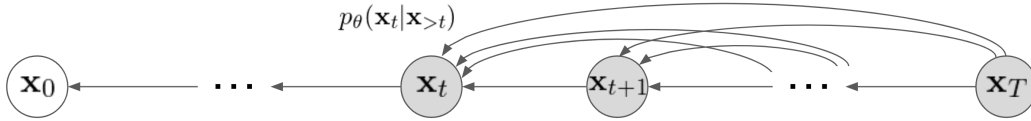


Figure 2: Illustration of GGDM. To improve sample quality, our novel family of samplers combines information from all previous (noisier) images at every denoising step.

- **VARS**: we parametrize the marginal variances of a DDPM as $\text{cumsum}(\text{softmax}([v; 1]))_t$ instead of fixing them to $1 - \bar{\alpha}_t$. This ensures they are monotonically increasing with respect to t (appending a one to ensure K degrees of freedom).
- **GGDM**: a new family of non-Markovian samplers for diffusion models with more degrees of freedom illustrated in Figure 2 and defined in Section 4. We parametrize μ_{tu} and σ_t of a GGDM for all t as sigmoid functions of trainable variables.
- **GGDM + PRED**: we parametrize all the μ_{tu} and σ_t identically to GGDM, but also learn the marginal coefficients with a $\text{cumsum} \circ \text{softmax}$ parametrization (identical to VARS) instead of computing them via Theorem 1, as well as the coefficients that predict \mathbf{x}_0 from $a_t \mathbf{x}_t - b_t \epsilon$ with $1 + \text{softplus}$ and softplus parametrizations.
- **[family]+TIME**: for any sampler family, we additionally parametrize the timesteps used to query the score model with a $\text{cumsum} \circ \text{softmax}$ parametrization (identical to VARS).

As we will show in the experiments, despite the fact that our pre-trained DDPMs are trained with discrete timesteps, learning the timesteps is still helpful. In principle, this should only be possible for DDPMs trained with continuous time sampling (Chen et al., 2021a; Song et al., 2021b; Kingma et al., 2021), but in practice we find that DDPMs trained with continuously embedded discrete timesteps are still well-behaved when applied at timesteps not present during training. We think this is due to the regularity of the sinusoidal positional encodings Vaswani et al. (2017) used in these model architectures and training with a sufficiently large number of timesteps T .

3.1 DIFFERENTIABLE SAMPLE QUALITY SCORES

We can differentiate through a stochastic sampler using the reparameterization trick, but the question of which objective to optimize still remains. Prior work has shown that optimizing log-likelihood can actually worsen sample quality and FID scores in the few-step regime (Watson et al., 2021; Song et al., 2021a). Thus, we instead design a *perceptual* loss which simply compares mean statistics between model samples and real samples in the neural network feature space. These types of objectives have been shown in prior work to better correlate with human perception of sample quality (Johnson et al., 2016; Heusel et al., 2017), which we also confirm in our experiments.

We rely on the representations of the penultimate layer of a pre-trained InceptionV3 classifier (Szegedy et al., 2016) and optimize the Kernel Inception Distance (KID) (Bińkowski et al., 2018). Let $\phi(\mathbf{x})$ denote the inception features of an image \mathbf{x} and p_ψ represent a diffusion sampler with trainable parameters ψ . For a linear kernel, which works best in our experiments, the objective is:

$$\mathcal{L}(\psi) = \mathbb{E}_{\mathbf{x}_p \sim p_\psi} \mathbb{E}_{\mathbf{x}'_p \sim p_\psi} \phi(\mathbf{x}_p)^\top \phi(\mathbf{x}'_p) - \mathbb{E}_{\mathbf{x}_p \sim p_\psi} \mathbb{E}_{\mathbf{x}_q \sim q} \phi(\mathbf{x}_p)^\top \phi(\mathbf{x}_q) \quad (7)$$

More generally, KID can be expressed as:

$$\mathcal{L}_{\text{KID}}(\psi) = \left\| \mathbb{E}_{\mathbf{x}_p \sim p_\psi} f^*(\mathbf{x}_p) - \mathbb{E}_{\mathbf{x}_q \sim q} f^*(\mathbf{x}_q) \right\|_2^2, \quad (8)$$

where $f^*(\mathbf{x}) = \mathbb{E}_{\mathbf{x}'_p \sim p_\psi} k_\phi(\mathbf{x}, \mathbf{x}'_p) - \mathbb{E}_{\mathbf{x}'_q \sim q} k_\phi(\mathbf{x}, \mathbf{x}'_q)$ is the witness function for any differentiable, positive definite kernel k , and $k_\phi(\mathbf{x}, \mathbf{y}) = k(\phi(\mathbf{x}), \phi(\mathbf{y}))$. Note that f^* attains the supremum of the MMD. To enable stochastic gradient descent, we use an unbiased estimator of KID using a minibatch of n model samples $\mathbf{x}_p^{(1)} \dots \mathbf{x}_p^{(n)} \sim p_\psi$ and n real samples $\mathbf{x}_q^{(1)} \dots \mathbf{x}_q^{(n)} \sim q$:

$$\frac{1}{n(n-1)} \sum_{i \neq j}^n k_\phi(\mathbf{x}_p^{(i)}, \mathbf{x}_p^{(j)}) - \frac{2}{n^2} \sum_{i=1}^n \sum_{j=1}^n k_\phi(\mathbf{x}_p^{(i)}, \mathbf{x}_q^{(j)}) + c, \quad (9)$$

where c is constant in ψ . Since the sampling chain of any Gaussian diffusion process admits using the reparametrization trick, our loss function is fully differentiable with respect to the trainable variables ψ . We empirically find that using the perceptual features is crucial; i.e., by trying $\phi(\mathbf{x}) = \mathbf{x}$ to compare images directly on pixel space rather than neural network feature space (as above), we observe that our method makes samples consistently worsen in apparent quality as training progresses.

3.2 GRADIENT REMATERIALIZATION

In order for backpropagation to be feasible under reasonable memory constraints, one final problem must be addressed: since we are taking gradients with respect to model samples, the cost in memory to maintain the state of the forward pass scales linearly with the number of inference steps, which can quickly become unfeasible considering the large size of typical DDPM architectures. To address this issue, we use gradient rematerialization (Kumar et al., 2019b). Instead of storing a particular computation’s output from the forward pass required by the backward pass computation, we recompute it on the fly. To trade $\mathcal{O}(K)$ memory cost for $\mathcal{O}(K)$ computation time, we simply rematerialize calls to the pre-trained DDPM (i.e., the estimated score function), but keep in memory all the progressively denoised images from the sampling chain. In JAX (Bradbury et al., 2018), this is trivial to implement by simply wrapping the score function calls with `jax.remat`.

4 GENERALIZED GAUSSIAN DIFFUSION MODELS

We now present Generalized Gaussian Diffusion Model (GGDM), our novel family of Gaussian diffusion processes that includes DDIM as a special case mentioned in section 3. We define a joint distribution with no independence assumptions

$$q_{\mu,\sigma}(\mathbf{x}_0, \dots, \mathbf{x}_T) = q(\mathbf{x}_0)q(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=1}^{T-1} q_{\mu,\sigma}(\mathbf{x}_t|\mathbf{x}_{>t}, \mathbf{x}_0) \quad (10)$$

where the new factors are defined as

$$q_{\mu,\sigma}(\mathbf{x}_t|\mathbf{x}_{>t}, \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t \left| \sum_{u \in S_t} \mu_{tu} \mathbf{x}_u, \sigma_t^2 \mathbf{I}_d \right.\right) \quad (11)$$

(letting $S_t = \{0, \dots, T\} \setminus \{1, \dots, t\}$ for notation compactness), with σ_t and μ_{tu} free parameters $\forall t \in \{1, \dots, T\}, u \in S_t$. In other words, when predicting the next, less noisy image, the sampler can take into account *all* the previous, noisier images in the sampling chain, and similarly to DDIM, we can also control the sampler’s variances. As we prove in the appendix (A.2), this construction admits Gaussian marginals, and we can differentially compute the marginal coefficients from arbitrary choices of μ and σ :

Theorem 1. Given some $t \in \{1, \dots, T\}$, let $a_{tu}^{(1)} = \mu_{tu} \forall u \in S_t$ and $v_t^{(1)} = \sigma_t^2$. For each $i \in \{1, \dots, T-t\}$, recursively define

$$a_{tu}^{(i+1)} = a_{t,t+i}^{(i)} \mu_{t+i,u} + a_{tu}^{(i)} \forall u \in S_{t+i} \quad \text{and} \quad v_t^{(i+1)} = v_t^{(i)} + \left(a_{t,t+i}^{(i)} \sigma_{t+i}\right)^2.$$

Then, it follows that

$$q_{\mu,\sigma}(\mathbf{x}_t|\mathbf{x}_{>t+i}, \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t \left| \sum_{u \in S_{t+i}} a_{tu}^{(i+1)} \mathbf{x}_u, v_t^{(i+1)} \mathbf{I}_d \right.\right). \quad (12)$$

In other words, instead of letting the β_t (or equivalently, the $\bar{\alpha}_t$) define the forward process as done by a usual DDPM, the GGDM family lets the μ_{tu} and σ_t define the process. In particular, an immediate corollary of Theorem 1 is that the marginal coefficients are given by

$$q_{\mu,\sigma}(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t \left| a_{t0}^{(T-t+1)} \mathbf{x}_0, v_t^{(T-t+1)} \mathbf{I}_d \right.\right) \quad (13)$$

The reverse process is thus defined as $p(\mathbf{x}_T) \prod_{t=1}^T p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ with $p(\mathbf{x}_T) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and

$$p_{\theta}(\mathbf{x}_t|\mathbf{x}_{>t}) = q_{\mu,\sigma}\left(\mathbf{x}_t|\mathbf{x}_{>t}, \hat{\mathbf{x}}_0 = \frac{1}{a_{t0}^{(T-t+1)}} \left(\mathbf{x}_t - \sqrt{v_t^{(T-t+1)}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)\right)\right). \quad (14)$$

Table 1: FID and IS scores for DDSS against baseline methods for a DDPM trained on CIFAR10 with the L_{simple} objective proposed by (Ho et al., 2020). FID scores (lower is better) are the numbers at the left of each entry, and IS scores (higher is better) are at the right.

Sampler \ K	5	10	15	20	25
DDPM (linear stride)	84.27 / 5.396	43.39 / 7.034	31.40 / 7.609	25.94 / 7.879	22.60 / 8.043
DDPM (quadratic stride)	76.25 / 5.435	42.03 / 6.965	27.78 / 7.714	20.225 / 8.128	16.17 / 8.350
DDIM (linear stride)	44.41 / 6.750	19.11 / 7.965	14.06 / 8.190	11.82 / 8.420	10.52 / 8.512
DDIM (quadratic stride)	32.66 / 7.090	13.62 / 8.190	9.318 / 8.495	7.500 / 8.641	6.560 / 8.759
GGDM +PRED+TIME	13.77 / 8.520	8.227 / 8.903	6.115 / 9.050	4.722 / 9.261	4.250 / 9.186

Table 2: FID / IS scores for DDSS against baseline methods for a DDPM trained on ImageNet 64x64 with the L_{hybrid} objective proposed by Nichol & Dhariwal (2021).

Sampler \ K	5	10	15	20	25
DDPM (linear stride)	122.0 / 5.878	58.78 / 10.67	39.30 / 13.22	31.36 / 14.72	26.36 / 15.71
DDPM (quadratic stride)	394.8 / 1.351	129.5 / 5.997	80.10 / 9.595	61.34 / 11.60	49.60 / 13.01
DDIM (linear stride)	135.4 / 5.898	40.70 / 12.225	28.54 / 13.99	24.225 / 14.75	22.13 / 15.16
DDIM (quadratic stride)	409.1 / 1.380	148.6 / 5.533	67.65 / 9.842	45.60 / 11.99	36.11 / 13.225
GGDM +PRED+TIME	55.14 / 12.90	37.32 / 14.76	24.69 / 17.225	20.69 / 17.92	18.40 / 18.12

4.1 IGNORING THE MATCHING MARGINALS CONDITION

Unlike DDIM, the GGDM family does not guarantee that the marginals of the new forward process match that of the original DDPM. We empirically find, however, that this condition can often be too restrictive and better samplers exist where the marginals $q_{\mu,\sigma}(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | a_{i_0}^{(T-t+1)} \mathbf{x}_0, v_t^{(T-t+1)} \mathbf{I}_d)$ of the new forward process differ from the original DDPM’s marginals. We verify this empirically by applying DDSS to both the family of DDIM sigmas and DDPM variances (“VARS” in Section 3): both sampler families have the same number of parameters (the reverse process variances), but the latter does not adjust the mean coefficients like DDIM to ensure matching marginals and still achieves similar or better scores than the former across sample quality metrics (and even outperforms the DDIM($\eta = 0$) baseline); see Section 5.2.

5 EXPERIMENTS

In order to emphasize that our method is compatible with any pre-trained DDPM, we apply our method on pre-trained DDPM checkpoints from prior work. Specifically, we experiment with the DDPM trained by Ho et al. (2020) with L_{simple} on CIFAR10, as well as a DDPM following the exact configuration of Nichol & Dhariwal (2021) trained on ImageNet 64x64 (Deng et al., 2009) with their L_{hybrid} objective (with the only difference being that we trained the latter ourselves for 3M rather than 1.5M steps). Both of these models utilize adaptations of the UNet architecture (Ronneberger et al., 2015) that incorporate self-attention layers (Vaswani et al., 2017).

We evaluate all of our models on both FID and Inception Score (IS) (Salimans et al., 2016), comparing the samplers discovered by DDSS against DDPM and DDIM baselines with linear and quadratic strides. As previously mentioned, more recent methods for fast sampling are outperformed by DDIM when the budget of inference steps is as small as those we utilize in this work (5, 10, 15, 20, 25). All reported results on both of these approximate sample quality metrics were computed by comparing 50K model and training data samples, as is standard in the literature. Also as is standard, IS scores are computed 10 times, each time on 5K samples, and then averaged.

In all of our experiments, we optimize the DDSS objective presented in Section 3.1 with the following hyperparameters:

1. For every family of models we search over, we initialize the degrees of freedom such that training begins with a sampler matching DDPM with K substeps following Song et al. (2020); Nichol & Dhariwal (2021).

Table 3: FID / IS scores for the KID kernel ablation on CIFAR10. When not learning the timesteps, we fix them to a quadratic stride, as Table 1 shows this performs best on CIFAR10.

Sampler \ K	5	10	15	20	25
DDSS (linear kernel)					
GGDM +PRED+TIME	13.77 / 8.520	8.227 / 8.903	6.115 / 9.050	4.722 / 9.261	4.250 / 9.186
GGDM +PRED	14.26 / 8.406	8.617 / 8.842	5.939 / 9.035	4.893 / 9.153	4.574 / 9.145
GGDM +TIME	12.85 / 8.383	7.858 / 8.895	6.265 / 9.075	5.367 / 9.136	4.887 / 9.229
GGDM)	14.45 / 8.281	8.154 / 8.892	7.045 / 8.939	5.477 / 9.183	4.815 / 9.189
DDSS (cubic kernel)					
GGDM +PRED+TIME	14.41 / 8.527	8.2257 / 9.007	5.895 / 9.036	4.932 / 9.092	4.278 / 9.286
GGDM +PRED	14.39 / 8.401	8.977 / 8.870	6.517 / 8.970	4.915 / 9.132	4.471 / 9.247
GGDM +TIME	12.35 / 8.406	7.879 / 8.852	6.682 / 8.999	5.639 / 9.058	4.631 / 9.189
GGDM	14.57 / 8.297	8.2252 / 8.836	6.727 / 8.904	5.569 / 9.177	4.668 / 9.192

- We apply gradient updates using the Adam optimizer (Kingma & Ba, 2015). We swept over the learning rate and used $\lambda = 0.0005$. We did not sweep over other Adam hyperparameters and kept $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 1 \times 10^{-8}$.
- We tried batch sizes of 128 and 512 and opted for the latter, finding that it leads to better sample quality upon inspection. Since the loss depends on averages over examples as our experiments are on unconditional generation, this choice was expected.
- We run all of our experiments for 50K training steps and evaluate the discovered samplers at this exact number of training steps. We did not sweep over this value.

We include our main results in Table 1 for CIFAR10 and Table 2 for ImageNet 64x64, comparing DDSS applied to GGDM +PRED+TIME against DDPM and DDIM baselines with linear and quadratic strides. All models use a linear kernel, i.e., $k_\phi(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^\top \phi(\mathbf{y})$, which we found to perform slightly better than the cubic kernel used by Bińkowski et al. (2018) (we ablate this in section 5.1). We omit the use of the learned variances of the ImageNet 64x64 model (i.e., following Nichol & Dhariwal (2021)), as we search for the variances ourselves via DDSS. We include samples for 5, 10 and 25 steps comparing the strongest DDIM baselines to DDSS + GGDM with a learned stride; see Figures 1 and 3. We include additional ImageNet 64x64 samples (A.1) and results for larger resolution datasets (A.4) in the appendix.

5.1 ABLATIONS FOR KID KERNEL AND GGDM VARIANTS

As our approach is compatible with any choice of KID kernel, we experiment with different choices of kernels. Namely, we try the simplest possible linear kernel, $k_\phi(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^\top \phi(\mathbf{y})$, as well as the cubic kernel $k_\phi(\mathbf{x}, \mathbf{y}) = \left(\frac{1}{d}\phi(\mathbf{x})^\top \phi(\mathbf{y}) + 1\right)^3$ used by Bińkowski et al. (2018). We compare the performance of these two kernels, as well as different variations of GGDM (i.e., with and without TIME and PRED as defined in Section 3). Results are included for CIFAR10 across all budgets in Table 3. We also include a smaller version of this ablation for ImageNet 64x64 in the appendix (A.3).

The results in this ablation show that the contributions of the linear kernel, timestep learning, and the empirical choice of learning the coefficients that predict \mathbf{x}_0 all slightly contribute to better FID and IS scores. Importantly, however, removing any of these additions still allows us to comfortably outperform the strongest baselines. See also the results on LSUN in the appendix A.4, which also do not include these additional trainable variables.

5.2 SEARCH SPACE ABLATION

Now, in order to further demonstrate the key importance of optimizing our GGDM family to find high-fidelity samplers, we also apply DDSS to the less general DDIM and VARS families. We show that, while we still attain better scores than a regular DDPM, searching these less flexible families of samplers does not yield improvements as significant as with our novel GGDM family. In particular, optimizing the DDIM sigma coefficients does not outperform the corresponding DDIM($\eta = 0$) baseline on CIFAR10, which is not a surprising result as Song et al. (2020) show empirically that most choices of the σ_t degrees of freedom lead to worse FID scores than setting them all to 0. These

Table 4: FID / IS scores for the DDSS search space ablation on CIFAR10. All runs fix the timesteps to a quadratic stride and use a linear kernel except for the last row (we only include the GGDM results for ease of comparison).

Sampler \ K	5	10	15	20	25
DDIM($\eta = 0$)	32.66 / 7.090	13.62 / 8.190	9.318 / 8.495	7.500 / 8.641	6.560 / 8.759
DDSS					
VARS	33.08 / 7.096	15.33 / 8.559	9.693 / 8.845	7.297 / 8.924	6.172 / 9.057
DDIM	32.61 / 7.084	16.29 / 7.966	11.31 / 8.372	9.120 / 8.563	7.853 / 8.644
GGDM	14.45 / 8.281	8.154 / 8.892	7.045 / 8.939	5.477 / 9.183	4.815 / 9.189
GGDM +PRED+TIME	13.77 / 8.520	8.227 / 8.903	6.115 / 9.050	4.722 / 9.261	4.250 / 9.186

results also show that optimizing the VARS can outperform DDSS applied to the DDIM family, and even the strongest DDIM($\eta = 0$) baselines for certain budgets, justifying our choice of not enforcing the marginals to match (as discussed in Section 4.1).

6 DISCUSSION

When applied to a sufficiently flexible family (such as the GGDM family proposed in this work), DDSS consistently finds samplers that achieve better image generation quality than the strongest baselines in the literature for very few steps. This is qualitatively apparent in non-cherry-picked samples (*e.g.*, DDIM($\eta = 0$) tends to generate blurrier images and with less background details as the budget decreases), and multiple quantitative sample quality metrics (FID and IS) also reflect these results. Still, we observe limitations to our method. Finding samplers with inference step budgets as small as $K < 10$ that have little apparent loss in quality remains challenging with our proposed search family. And, while on CIFAR10 the metrics indicate significant relative improvement over sample quality metrics, the relative improvement on ImageNet 64x64 is less pronounced. We hypothesize that this is an inherent difficulty of ImageNet due to its high diversity of samples, and that in order to retain sample quality and diversity, it might be impossible to escape some minimum number of inference steps with score-based models as they might be crucial to mode-breaking.

Beyond the empirical gains of applying our procedure, our findings shed further light into properties of pre-trained score-based generative models. First, we show that without fine-tuning a DDPM’s parameters, these models are already capable of producing high-quality samples with very few inference steps, though the default DDPM sampler in this regime is usually suboptimal when using a few-step sampler. We further show that better sampling paths exist, and interestingly, these are determined by alternative variational lower bounds to the data distribution that make use of the score-based model but do not necessarily share the same marginals as the original DDPM forward process. Our findings thus suggest that enforcing this marginal-sharing constraint is unnecessary and can be too restrictive in practice.

7 OTHER RELATED WORK

Besides DDIM (Song et al., 2020), there have been more recent attempts at reducing the number of inference steps for DDPMs. Jolicœur-Martineau et al. (2021) proposed a dynamic step size SDE solver that can reduce the number of calls to the modeled score function to ~ 150 on CIFAR10 (Krizhevsky et al., 2009) with minimal cost in FID scores, but quickly falls behind DDIM($\eta = 0$) with as many as 50 steps. Watson et al. (2021) proposed a dynamic programming algorithm that chooses log-likelihood optimal strides, but find that log-likelihood reduction has a mismatch with FID scores, particularly with in the very few step regime, also falling behind DDIM($\eta = 0$) in this front. Other methods that have been shown to help sample quality in the few-step regime include non-Gaussian variants of diffusion models (Nachmani et al., 2021) and adaptively adjusting the sampling path by introducing a noise level estimating network (San-Roman et al., 2021), but more thorough evaluation of sample quality achieved by these approaches is needed with budgets as small as those considered in this work.

Other approaches to sampling DDPMs have also been recently proposed, though not for the explicit purpose of efficient sampling. Song et al. (2021b) derive a reverse SDE that, when discretized, uses

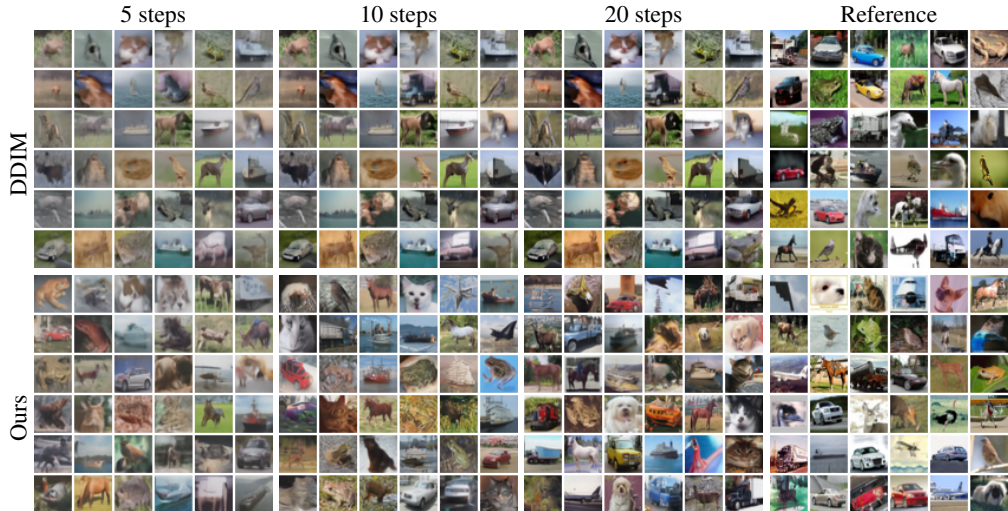


Figure 3: Non-cherrypicked samples for a DDPM trained on CIFAR10, comparing the strongest DDIM($\eta = 0$) baseline and our approach. All samples were generated with the same random seeds. For reference, we include DDPM samples using all 1,000 steps (top right) and real images (bottom right).

different coefficients than the ancestral samplers considered in this work. The same authors also derive “corrector” steps, which introduce additional calls to the pre-trained DDPM as a form of gradient ascent (Langevin dynamics) that help with quality but introduce computation cost, as well as an alternative sampling procedure using a probability flow ODE that shares the same marginals as the DDPM’s original forward process. Huang et al. (2021) generalize this family of samplers to a “plug-in reverse SDE” that interpolates between a probability flow ODE and the reverse SDE, similarly to how the DDIM η interpolates between an implicit probabilistic model and a stochastic reverse process. Our proposed search family includes discretizations of most of these cases for Gaussian processes, notably missing corrector steps, where reusing a single timestep is considered.

8 CONCLUSION AND FUTURE WORK

We propose Differentiable Diffusion Sampler Search (DDSS), a method for finding few-step samplers for Denoising Diffusion Probabilistic Models. We show how to optimize a perceptual loss over a space of diffusion processes that makes use of a pre-trained DDPM’s samples by leveraging the reparametrization trick and gradient rematerialization. Our results qualitatively and quantitatively show that DDSS is able to significantly improve sample quality for unconditional image generation over prior methods on efficient DDPM sampling. The success of our method hinges on searching a novel, wider family of Generalized Gaussian Diffusion Model (GGDM) than those identified in prior work (Song et al., 2020). DDSS does not fine-tune the pre-trained DDPM, only needs to be applied once, has few hyperparameters, and does not require re-training the DDPM.

Our findings suggest future directions to further reduce the number of inference steps while retaining high fidelity in generated samples. For instance, it is plausible to use different representations for the perceptual loss instead of those of a classifier, *e.g.*, use representations from an unsupervised model such as SimCLR (Chen et al., 2020), to using internal representations learned by the pre-trained DDPM itself, which would eliminate the burden of additional computation. Moreover, considering the demonstrated benefits of applying DDSS to our proposed GGDM family of samplers (as opposed to narrower families like DDIM), we motivate future work on identifying more general families of samplers and investigating whether they help uncover even better samplers or lead to overfitting. Finally, identifying other variants of perceptual losses (*e.g.*, that do not sample from the model), or alternative optimization strategies (*e.g.*, gradient-free methods) that lead to similar results is important future work. This would make DDSS itself a more efficient procedure, as gradient-based optimization of our proposed loss requires extensive memory or computation requirements to back-propagate through the whole sampling chain.

REFERENCES

- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. In *arXiv*, 2017.
- Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. *arXiv preprint arXiv:1801.01401*, 2018.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- Ruojin Cai, Guandao Yang, Hadar Averbuch-Elor, Zekun Hao, Serge Belongie, Noah Snaveley, and Bharath Hariharan. Learning Gradient Fields for Shape Generation. In *ECCV*, 2020.
- Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan. WaveGrad: Estimating Gradients for Waveform Generation. In *ICLR*, 2021a.
- Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, Najim Dehak, and William Chan. WaveGrad 2: Iterative Refinement for Text-to-Speech Synthesis. In *INTERSPEECH*, 2021b.
- Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big Self-Supervised Models are Strong Semi-Supervised Learners. In *NeurIPS*, 2020.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. *arXiv preprint arXiv:2105.05233*, 2021.
- Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of wasserstein gans. *arXiv preprint arXiv:1704.00028*, 2017.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *arXiv preprint arXiv:1706.08500*, 2017.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. *NeurIPS*, 2020.
- Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *arXiv preprint arXiv:2106.15282*, 2021.
- Chin-Wei Huang, Jae Hyun Lim, and Aaron Courville. A variational perspective on diffusion-based generative models and score matching. *arXiv preprint arXiv:2106.02808*, 2021.
- Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pp. 694–711. Springer, 2016.
- Alexia Jolicoeur-Martineau, Ke Li, Rémi Piché-Taillefer, Tal Kachman, and Ioannis Mitliagkas. Gotta go fast when generating data with score-based models. *arXiv preprint arXiv:2105.14080*, 2021.
- Diederik Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *ICLR*, 2015.

- Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. In *ICLR*, 2013.
- Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *arXiv preprint arXiv:2107.00630*, 2021.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Technical Report*, 2009.
- Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre de Brebisson, Yoshua Bengio, and Aaron Courville. MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis. In *NeurIPS*, 2019a.
- Ravi Kumar, Manish Purohit, Zoya Svitkina, Erik Vee, and Joshua R Wang. Efficient rematerialization for deep networks. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pp. 15172–15181, 2019b.
- Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2837–2845, 2021.
- Eliya Nachmani, Robin San Roman, and Lior Wolf. Non gaussian denoising diffusion models. *arXiv preprint arXiv:2106.07582*, 2021.
- Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *arXiv preprint arXiv:2102.09672*, 2021.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29: 2234–2242, 2016.
- Robin San-Roman, Eliya Nachmani, and Lior Wolf. Noise estimation for generative diffusion models. *arXiv preprint arXiv:2104.02600*, 2021.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pp. 2256–2265. PMLR, 2015.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- Yang Song and Stefano Ermon. Generative Modeling by Estimating Gradients of the Data Distribution. *NeurIPS*, 2019.
- Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *arXiv e-prints*, pp. arXiv–2101, 2021a.
- Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations. In *ICLR*, 2021b.
- Markus Svensén and Christopher M Bishop. Pattern recognition and machine learning, 2007.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Re-thinking the Inception Architecture for Computer Vision. In *CVPR*, 2016.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *NIPS*, 2017.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011.

Daniel Watson, Jonathan Ho, Mohammad Norouzi, and William Chan. Learning to efficiently sample from diffusion probabilistic models. *arXiv preprint arXiv:2106.03802*, 2021.

Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.

A APPENDIX

A.1 ADDITIONAL IMAGENET 64X64 SAMPLES

We provide additional samples for our results on ImageNet 64x64. The DDPM and DDIM($\eta = 0$) samples (left and middle, respectively) use a linear stride, while our DDSS + GGDM samples (right) use a learned stride.



Figure A.1: Additional samples on ImageNet 64x64. For reference, we include DDPM samples with all 4,000 steps (bottom left) and real samples (bottom middle).

A.2 PROOF FOR THEOREM 1

Theorem 1. Given some $t \in \{1, \dots, T\}$, let $a_{tu}^{(1)} = \mu_{tu} \forall u \in S_t$ and $v_t^{(1)} = \sigma_t^2$. For each $i \in \{1, \dots, T-t\}$, recursively define

- $a_{tu}^{(i+1)} = a_{t,t+i}^{(i)} \mu_{t+i,u} + a_{tu}^{(i)} \forall u \in S_{t+i}$
- $v_t^{(i+1)} = v_t^{(i)} + \left(a_{t,t+i}^{(i)} \sigma_{t+i} \right)^2$

Then, it follows that

$$q_{\mu,\sigma}(\mathbf{x}_t | \mathbf{x}_{>t+i}, \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_{t+i}} a_{tu}^{(i+1)} \mathbf{x}_u, v_t^{(i+1)} \mathbf{I}_d \right. \right)$$

Proof. Let us prove this result with mathematical induction. Note that, for each such t , we have by definition that

$$q_{\mu,\sigma}(\mathbf{x}_{t+1} | \mathbf{x}_{>t+1}, \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_{t+1}} \mu_{t+1,u} \mathbf{x}_u, \sigma_{t+1}^2 \mathbf{I}_d \right. \right)$$

and

$$q_{\mu,\sigma}(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{x}_{>t+1}, \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_t} \mu_{tu} \mathbf{x}_u, \sigma_t^2 \mathbf{I}_d \right. \right) = \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_t} a_{tu}^{(1)} \mathbf{x}_u, v_t^{(1)} \mathbf{I}_d \right. \right)$$

Therefore, following Svensén & Bishop (2007) (2.115), by prior conjugacy it follows that

$$\begin{aligned} q_{\mu,\sigma}(\mathbf{x}_t | \mathbf{x}_{>t+1}, \mathbf{x}_0) &= \mathcal{N} \left(\mathbf{x}_t \left| a_{t,t+1}^{(1)} \sum_{u \in S_{t+1}} \mu_{t+1,u} \mathbf{x}_u + \sum_{u \in S_{t+1}} a_{tu}^{(1)}, \left(v_t^{(1)} + a_{t,t+1}^{(1)} \sigma_{t+1}^2 a_{t,t+1}^{(1)} \right) \mathbf{I}_d \right. \right) \\ &= \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_{t+1}} \left(a_{t,t+1}^{(1)} \mu_{t+1,u} + a_{tu}^{(1)} \right) \mathbf{x}_u, v_t^{(2)} \mathbf{I}_d \right. \right) \\ &= \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_{t+1}} a_{tu}^{(2)} \mathbf{x}_u, v_t^{(2)} \mathbf{I}_d \right. \right). \end{aligned}$$

This proves the base case for our induction argument. Now, let us prove the inductive step. Suppose there exists some integer $j \in \{1, \dots, T-t+1\}$ such that

$$q_{\mu,\sigma}(\mathbf{x}_t | \mathbf{x}_{>t+j}, \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_{t+j}} a_{tu}^{(j+1)} \mathbf{x}_u, v_t^{(j+1)} \mathbf{I}_d \right. \right).$$

By definition, we already know $q(\mathbf{x}_{t+j+1} | \mathbf{x}_{>t+j+1}, \mathbf{x}_0)$, so we have

$$q_{\mu,\sigma}(\mathbf{x}_{t+j+1} | \mathbf{x}_{>t+j+1}, \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_{t+j+1} \left| \sum_{u \in S_{t+j+1}} \mu_{t+j+1,u} \mathbf{x}_u, \sigma_{t+j+1}^2 \mathbf{I}_d \right. \right)$$

and (rewriting the inductive hypothesis)

$$q_{\mu,\sigma}(\mathbf{x}_t | \mathbf{x}_{t+j+1}, \mathbf{x}_{>t+j+1}, \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_{t+j}} a_{tu}^{(j+1)} \mathbf{x}_u, v_t^{(j+1)} \mathbf{I}_d \right. \right).$$

Therefore, by prior conjugacy again, it follows that

$$\begin{aligned} & q_{\mu,\sigma}(\mathbf{x}_t | \mathbf{x}_{>t+j+1}, \mathbf{x}_0) \\ &= \mathcal{N} \left(\mathbf{x}_t \left| a_{t,t+j+1}^{(j+1)} \sum_{u \in S_{t+j}} \mu_{t+j+1,u} \mathbf{x}_u + \sum_{u \in S_{t+j}} a_{tu}^{(j+1)}, \left(v_t^{(j+1)} + a_{t,t+j+1}^{(j+1)} \sigma_{t+j+1}^2 a_{t,t+j+1}^{(j+1)} \right) \mathbf{I}_d \right. \right) \\ &= \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_{t+j}} \left(a_{t,t+j+1}^{(j+1)} \mu_{t+j+1,u} + a_{tu}^{(j+1)} \right) \mathbf{x}_u, v_t^{(j+2)} \mathbf{I}_d \right. \right) \\ &= \mathcal{N} \left(\mathbf{x}_t \left| \sum_{u \in S_{t+j+1}} a_{tu}^{(j+2)} \mathbf{x}_u, v_t^{(j+2)} \mathbf{I}_d \right. \right). \end{aligned}$$

This concludes the proof of the inductive step. Hence, we have proven the result for any $i \in \{1, \dots, T-t\}$. In particular,

$$q_{\mu,\sigma}(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_t \left| a_{t0}^{(T-t+1)} \mathbf{x}_0, v_t^{(T-t+1)} \mathbf{I}_d \right. \right). \quad \square$$

A.3 ADDITIONAL ABLATION OF KID KERNEL AND GGDM VARIANTS FOR IMAGENET 64x64

We also ran a smaller version of the ablation results presented in Section 5.1, but for ImageNet 64x64 instead of CIFAR10, as these are more computationally intensive to do a full grid search. Results for a step budget $K = 15$ are included below. When not learning the timesteps, we fix them to a linear stride, as Table 2 shows this performs best on ImageNet 64x64.

Sampler \ K	15
DDSS (linear kernel)	
GGDM +PRED+TIME	24.69 / 17.225
GGDM +PRED	27.08 / 16.44
GGDM +TIME	25.73 / 17.27
GGDM	28.34 / 16.63
DDSS (cubic kernel)	
GGDM +PRED+TIME	26.52 / 16.29
GGDM +PRED	27.82 / 16.3
GGDM +TIME	26.87 / 16.99
GGDM	28.83 / 16.32

A.4 RESULTS ON LARGER RESOLUTION DATASETS

We include results for LSUN (Yu et al., 2015) bedrooms and churches at the 128x128 resolution. We trained the models for 400K and 200K steps (respectively), and all other hyperparameters are identical: we use the Adam optimizer with learning rate 0.0003 (linearly warmed up for the first 1000 training steps), batch size 2048, gradient clipping at norms over 1.0, dropout of 0.1, and EMA over the weights with decay rate 0.9999. We train the models using a linear stride of 1000 evenly-spaced timesteps, fixing the log-signal-to-noise-ratio schedule to a cosine function monotonically decreasing from 20 to -20. The ELBO is reweighted with L_{simple} following Ho et al. (2020), but we additionally reweight each term by $\max(1, \text{SNR})$ which we found to be slightly helpful in resulting FID scores (note this is equivalent to minimizing the worst mean squared error between either the \mathbf{x}_0 or ϵ). The UNet employs five down/up-sampling resolutions with $768 \times (1, 2, 4, 6, 8)$ respective channels, 3 ResNet blocks per resolution, and spatial self-attention at the 3 smallest resolutions, i.e., 8, 16, and 32.

After training the models, we run DDSS using just the GGDM model family for simplicity (i.e., we don't use the +PRED and +TIME we experiment with in the paper) at 5, 10 and 20 evenly-spaced inference steps. DDSS training occurs for 50K steps, using the Adam optimizer with learning rate of 0.0005 and batch size 512, optimizing the linear kernel for the KID loss. We compare against the usual DDPM and DDIM($\eta = 0$) baselines at the same inference budgets and include the FID scores with all 1000 steps for reference. Results and samples are included below.

Sampler \ K	5	10	20	1000
LSUN Bedroom				
DDPM	95.38	44.84	16.88	2.547
DDIM($\eta = 0$)	168.7	56.33	9.527	-
DDSS (GGDM)	29.15	11.01	4.817	-
LSUN Church				
DDPM	96.67	51.05	16.53	2.718
DDIM($\eta = 0$)	133.1	54.39	14.96	-
DDSS (GGDM)	30.24	11.59	6.736	-



Figure A.2: Non-cherrypicked samples for a DDPM trained on LSUN bedroom 128x128, comparing DDPM and DDIM($\eta = 0$) to our approach. All samples were generated with the same random seeds and a linear stride. For reference, we include DDPM samples using all 1,000 steps (bottom left) and real images (bottom middle).



Figure A.3: Non-cherrypicked samples for a DDPM trained on LSUN church 128x128, comparing DDPM and DDIM($\eta = 0$) to our approach. All samples were generated with the same random seeds and a linear stride. For reference, we include DDPM samples using all 1,000 steps (bottom left) and real images (bottom middle).