

# PROGRESSIVE DISTILLATION FOR FAST SAMPLING OF DIFFUSION MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Diffusion models have recently shown great promise for generative modeling, outperforming GANs on perceptual quality and autoregressive models at density estimation. A remaining downside is their slow sampling time: generating high quality samples takes many hundreds or thousands of model evaluations. Here we make two contributions to help eliminate this downside: First, we present new parameterizations of diffusion models that provide increased stability when using few sampling steps, compared to models in the literature. Second, we present a method to distill a trained deterministic diffusion sampler, using many steps, into a new diffusion model that takes half as many sampling steps. We then keep progressively applying this distillation procedure to our model, halving the number of required sampling steps each time. On standard image generation benchmarks like CIFAR-10, ImageNet, and LSUN, we start out with (near) state-of-the-art samplers taking 1024 or 8192 steps, and are able to distill down to models taking as little as 4 steps without losing much perceptual quality; achieving, for example, a FID of 3.0 on CIFAR-10 in 4 steps. Finally, we show that the full progressive distillation procedure does not take more time than it takes to train the original model, thus representing an efficient solution for generative modeling using diffusion at both train and test time.

## 1 INTRODUCTION

Diffusion models (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020) are a class of generative models that have recently delivered competitive sample quality and likelihoods on many standard generative modeling benchmarks. Diffusion models have achieved ImageNet generation results outperforming BigGAN-deep and VQ-VAE-2 in terms of FID score and classification accuracy score (Ho et al., 2021; Dhariwal & Nichol, 2021a), and they have achieved likelihoods outperforming autoregressive image models (Kingma et al., 2021; Song et al., 2021b). Diffusion models have also succeeded in image super-resolution (Saharia et al., 2021; Li et al., 2021) and image inpainting (Song et al., 2021c), and there have been promising results in shape generation (Cai et al., 2020), graph generation (Niu et al., 2020), and text generation (Hoogeboom et al., 2021; Austin et al., 2021).

A major barrier remains to practical adoption of diffusion models: sampling speed. While sampling can be accomplished in relatively few steps in strongly conditioned settings, such as text-to-speech (Chen et al., 2021) and image super-resolution (Saharia et al., 2021), the situation is substantially different in scenarios in which there is little mutual information between the generated data and the conditioning signal. Examples of such scenarios are unconditional and class-conditional image generation, and results for these scenarios so far in the diffusion literature require hundreds or thousands of sampling steps using network evaluations that are not amenable to caching optimizations of other types of generative models (Ramachandran et al., 2017).

In this paper, we reduce the sampling time of diffusion models by orders of magnitude in unconditional and class-conditional image generation, which represent the scenarios in which diffusion models have been slowest in previous work. We present a procedure to distill the behavior of a  $T$ -step DDIM sampler (Song et al., 2021a) for a pretrained diffusion model into a new model with  $T/2$  steps, with little degradation in sample quality. In what we call *progressive distillation*, we repeat

this distillation procedure to produce models that generate in as little as 4 steps, still maintaining sample quality competitive with state-of-the-art models using thousands of steps.

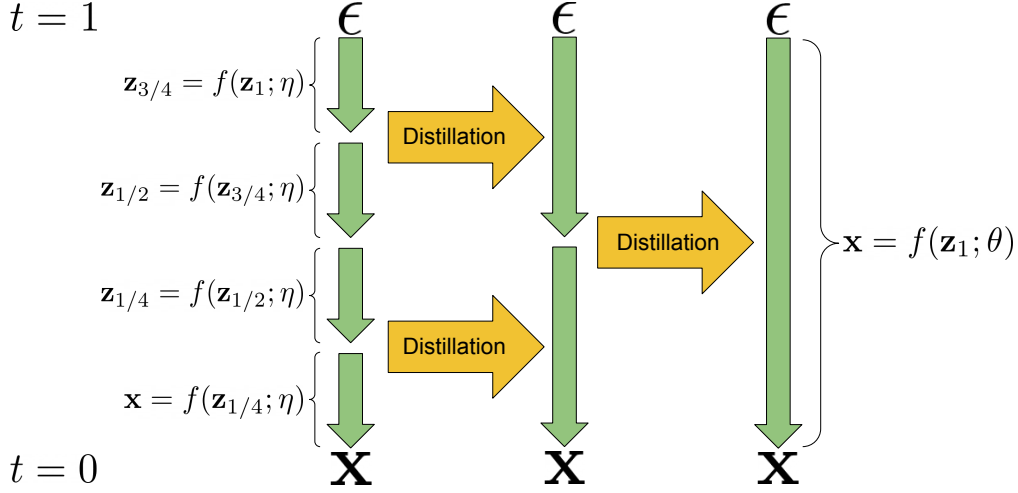


Figure 1: A visualization of performing two iterations of our proposed *progressive distillation* algorithm. A sampler  $f(\mathbf{z}; \eta)$ , mapping random noise  $\epsilon$  to samples  $\mathbf{x}$  in 4 deterministic steps, is distilled into a new sampler  $f(\mathbf{z}; \theta)$  taking only a single step. The sampler  $f(\mathbf{z}; \eta)$  is derived by approximately integrating the *probability flow ODE* for a learned diffusion model, and distillation can thus be understood as learning to integrate in fewer steps, or *amortizing* this integration into the sampler’s parameters  $\theta$ .

## 2 BACKGROUND ON DIFFUSION MODELS

We consider diffusion models (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020) specified in continuous time (Tzen & Raginsky, 2019a; Song et al., 2021c; Chen et al., 2021; Kingma et al., 2021). We use  $\mathbf{x} \sim p(\mathbf{x})$  to denote training data. A diffusion model has latent variables  $\mathbf{z} = \{\mathbf{z}_t | t \in [0, 1]\}$  and is specified by a noise schedule comprising differentiable functions  $\alpha_t, \sigma_t$  such that  $\lambda_t = \log[\alpha_t^2/\sigma_t^2]$ , the log signal-to-noise-ratio, decreases monotonically with  $t$ .

These ingredients define the forward process  $q(\mathbf{z}|\mathbf{x})$ , a Gaussian process satisfying the following Markovian structure:

$$q(\mathbf{z}_t|\mathbf{x}) = \mathcal{N}(\mathbf{z}_t; \alpha_t \mathbf{x}, \sigma_t^2 \mathbf{I}), \quad q(\mathbf{z}_t|\mathbf{z}_s) = \mathcal{N}(\mathbf{z}_t; (\alpha_t/\alpha_s)\mathbf{z}_s, \sigma_{t|s}^2 \mathbf{I}) \quad (1)$$

where  $0 \leq s < t \leq 1$  and  $\sigma_{t|s}^2 = (1 - e^{\lambda_t - \lambda_s})\sigma_t^2$ .

The role of function approximation in the diffusion model is to denoise  $\mathbf{z}_t \sim q(\mathbf{z}_t|\mathbf{x})$  into an estimate  $\hat{\mathbf{x}}_\theta(\mathbf{z}_t) \approx \mathbf{x}$  (the function approximator also receives  $\lambda_t$  as an input, but we omit this to keep our notation clean). We train this denoising model  $\hat{\mathbf{x}}_\theta$  using a weighted mean squared error loss

$$\mathbb{E}_{\epsilon, t} [w(\lambda_t) \|\hat{\mathbf{x}}_\theta(\mathbf{z}_t) - \mathbf{x}\|_2^2] \quad (2)$$

over uniformly sampled times  $t \in [0, 1]$ . This loss can be justified as a weighted variational lower bound on the data log likelihood under the diffusion model (Kingma et al., 2021) or as a form of denoising score matching (Vincent, 2011; Song & Ermon, 2019). We will discuss particular choices of weighting functions later on in the paper.

Sampling from a trained model can be performed in several ways. The most straightforward way is discrete time ancestral sampling (Ho et al., 2020). To define this sampler, first note that the forward process can be described in reverse as  $q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{x}) = \mathcal{N}(\mathbf{z}_s; \tilde{\boldsymbol{\mu}}_{s|t}(\mathbf{z}_t, \mathbf{x}), \tilde{\sigma}_{s|t}^2 \mathbf{I})$  (noting  $s < t$ ), where

$$\tilde{\boldsymbol{\mu}}_{s|t}(\mathbf{z}_t, \mathbf{x}) = e^{\lambda_t - \lambda_s} (\alpha_s/\alpha_t) \mathbf{z}_t + (1 - e^{\lambda_t - \lambda_s}) \alpha_s \mathbf{x}, \quad \tilde{\sigma}_{s|t}^2 = (1 - e^{\lambda_t - \lambda_s}) \sigma_s^2 \quad (3)$$

We use this reversed description of the forward process to define the ancestral sampler. Starting at  $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , the ancestral sampler follows the rule

$$\mathbf{z}_s = \tilde{\boldsymbol{\mu}}_{s|t}(\mathbf{z}_t, \hat{\mathbf{x}}_\theta(\mathbf{z}_t)) + \sqrt{(\tilde{\sigma}_{s|t}^2)^{1-v}(\sigma_{t|s}^2)^v} \boldsymbol{\epsilon} \quad (4)$$

$$= e^{\lambda_t - \lambda_s} (\alpha_s / \alpha_t) \mathbf{z}_t + (1 - e^{\lambda_t - \lambda_s}) \alpha_s \hat{\mathbf{x}}_\theta(\mathbf{z}_t) + \sqrt{(\tilde{\sigma}_{s|t}^2)^{1-v}(\sigma_{t|s}^2)^v} \boldsymbol{\epsilon} \quad (5)$$

where  $v$  is a hyperparameter that controls how much noise is added (Nichol & Dhariwal, 2021a).

Alternatively, Song et al. (2021c) show that our learned denoising model  $\hat{\mathbf{x}}_\theta(\mathbf{z}_t)$  can be used to deterministically map noise  $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  to samples  $\mathbf{x}$  by numerically solving the *probability flow ODE*:

$$d\mathbf{z}_t = [f(\mathbf{z}_t, t) - \frac{1}{2}g^2(t)\nabla_z \log \hat{p}_\theta(\mathbf{z}_t)]dt, \quad (6)$$

where  $\nabla_z \log \hat{p}_\theta(\mathbf{z}_t) = \frac{\alpha_t \hat{\mathbf{x}}_\theta(\mathbf{z}_t) - \mathbf{z}_t}{\sigma_t^2}$ . Following Kingma et al. (2021), we have  $f(\mathbf{z}_t, t) = \frac{d \log \alpha_t}{dt} \mathbf{z}_t$  and  $g^2(t) = \frac{d\sigma_t^2}{dt} - 2 \frac{d \log \alpha_t}{dt} \sigma_t^2$ . Since  $\hat{\mathbf{x}}_\theta(\mathbf{z}_t)$  is parameterized by a neural network, this equation is a special case of a *neural ODE* (Chen et al., 2018), also called a *continuous normalizing flow* (Grathwohl et al., 2018).

Solving the ODE in Equation 6 numerically can be done with standard methods like the Euler rule or the Runge-Kutta method. The DDIM sampler proposed by Song et al. (2021a) can also be understood as an integration rule for this ODE, as we show in Appendix B, even though it was originally proposed with a different motivation. The update rule specified by DDIM is

$$\mathbf{z}_s = \alpha_s \hat{\mathbf{x}}_\theta(\mathbf{z}_t) + \sigma_s \frac{\mathbf{z}_t - \alpha_t \hat{\mathbf{x}}_\theta(\mathbf{z}_t)}{\sigma_t} \quad (7)$$

$$= e^{(\lambda_t - \lambda_s)/2} (\alpha_s / \alpha_t) \mathbf{z}_t + (1 - e^{(\lambda_t - \lambda_s)/2}) \alpha_s \hat{\mathbf{x}}_\theta(\mathbf{z}_t), \quad (8)$$

and in practice this rule performs better than the aforementioned standard ODE integration rules in our case, as we show in Appendix C.

If  $\hat{\mathbf{x}}_\theta(\mathbf{z}_t)$  satisfies mild smoothness conditions, the error introduced by numerical integration of the probability flow ODE is guaranteed to vanish as the number of integration steps grows infinitely large, i.e.  $N \rightarrow \infty$ . This leads to a trade-off in practice between the accuracy of the numerical integration, and hence the quality of the produced samples from our model, and the time needed to produce these samples. So far, most models in the literature have needed hundreds or thousands of integration steps to produce their highest quality samples, which is prohibitive for many practical applications of generative modeling. Here, we therefore propose a method to distill these accurate, but slow, ODE integrators into much faster models that are still very accurate. This idea is visualized in Figure 1, and described in detail in the next section.

### 3 PROGRESSIVE DISTILLATION

To make diffusion models more efficient at sampling time, we propose *progressive distillation*: an algorithm that iteratively halves the number of required sampling steps by distilling a slow teacher diffusion model into a faster student model. Our implementation of progressive distillation stays very close to the implementation for training the original diffusion model, as described by e.g. Ho et al. (2020). Algorithm 1 and Algorithm 2 present diffusion model training and progressive distillation side-by-side, with the relative changes in progressive distillation highlighted in **green**.

We start the progressive distillation procedure with a teacher diffusion model that is obtained by training it in the standard way. At every iteration of progressive distillation we then initialize the student model with a copy of the teacher, using both the same parameters and same model definition. Like in standard training, we then sample data from the training set and add noise to it, before forming the training loss by applying the student denoising model to this noisy data  $\mathbf{z}_t$ . The main difference in progressive distillation is in how we set the target for the denoising model: instead of the original data  $\mathbf{x}$ , we have the student model denoise towards a target  $\tilde{\mathbf{x}}$  that makes a single student DDIM step match 2 teacher DDIM steps. We calculate this target value by running 2 DDIM sampling steps using the teacher, starting from  $\mathbf{z}_t$  and ending at  $\mathbf{z}_{t-1/N}$ , when distilling to  $N$  student

sampling steps. By inverting a single step of DDIM, we then calculate the value the student model would need to predict in order to step from  $\mathbf{z}_t$  to  $\mathbf{z}_{t-1/N}$  in a single step. The resulting target value  $\tilde{\mathbf{x}}(\mathbf{z}_t)$  is deterministically determined given the teacher model and starting point  $\mathbf{z}_t$ , unlike the original data point  $\mathbf{x}$ , since multiple different data points  $\mathbf{x}$  could conceivably have led to observing noisy data  $\mathbf{z}_t$ . This allows the student to make a sharper prediction than the teacher, thus making faster progress. After running distillation to learn a student model taking  $N$  sampling steps, we can repeat the procedure with  $N/2$  steps: The student model then becomes the new teacher, and a new student model is initialized by making a copy of it.

Unlike our procedure for training the original model, we always run progressive distillation in discrete time: we sample this discrete time such that the highest time index corresponds to a signal-to-noise ratio of zero, i.e.  $\alpha_1 = 0$ , which exactly matches the distribution of input noise  $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  that is used at test time. We found this to work slightly better than starting from a non-zero signal-to-noise ratio as used by e.g. Ho et al. (2020).

---

**Algorithm 1** Standard diffusion training

---

**Require:** Model  $\hat{\mathbf{x}}_\theta(\mathbf{z}_t)$  to be trained

**Require:** Data set  $\mathcal{D}$

**Require:** Loss weight function  $w()$

**while** not converged **do**

$\mathbf{x} \sim \mathcal{D}$   $\triangleright$  Sample data

$t \sim U[0, 1]$   $\triangleright$  Sample time

$\epsilon \sim N(0, I)$   $\triangleright$  Sample noise

$\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$   $\triangleright$  Add noise to data

$\tilde{\mathbf{x}} = \mathbf{x}$   $\triangleright$  Clean data is target for  $\hat{\mathbf{x}}$

$\lambda_t = \log[\alpha_t^2 / \sigma_t^2]$   $\triangleright$  log-SNR

$L_\theta = w(\lambda_t) \|\tilde{\mathbf{x}} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t)\|_2^2$   $\triangleright$  Loss

$\theta \leftarrow \theta - \gamma \nabla_\theta L_\theta$   $\triangleright$  Optimization

**end while**

---



---

**Algorithm 2** Progressive distillation

---

**Require:** Trained teacher model  $\hat{\mathbf{x}}_\eta(\mathbf{z}_t)$

**Require:** Data set  $\mathcal{D}$

**Require:** Loss weight function  $w()$

**Require:** Student sampling steps  $N$

**for**  $K$  iterations **do**

$\theta \leftarrow \eta$

$\triangleright$  Init student from teacher

**while** not converged **do**

$\mathbf{x} \sim \mathcal{D}$

$t = i/N, i \sim \text{Cat}[1, 2, \dots, N]$

$\epsilon \sim N(0, I)$

$\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$

**# 2 steps of DDIM with teacher**

$t' = t - 0.5/N, t'' = t - 1/N$

$\mathbf{z}_{t'} = \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_t) + \frac{\sigma_{t'}}{\sigma_t} (\mathbf{z}_t - \alpha_t \hat{\mathbf{x}}_\eta(\mathbf{z}_t))$

$\mathbf{z}_{t''} = \alpha_{t''} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}) + \frac{\sigma_{t''}}{\sigma_{t'}} (\mathbf{z}_{t'} - \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}))$

$\tilde{\mathbf{x}} = \frac{\mathbf{z}_{t''} - (\sigma_{t''}/\sigma_t) \mathbf{z}_t}{\alpha_{t''} - (\sigma_{t''}/\sigma_t) \alpha_t}$   $\triangleright$  Teacher  $\hat{\mathbf{x}}$  target

$\lambda_t = \log[\alpha_t^2 / \sigma_t^2]$

$L_\theta = w(\lambda_t) \|\tilde{\mathbf{x}} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t)\|_2^2$

$\theta \leftarrow \theta - \gamma \nabla_\theta L_\theta$

**end while**

$\eta \leftarrow \theta$

$\triangleright$  Student becomes next teacher

$N \leftarrow N/2$   $\triangleright$  Halve number of sampling steps

**end for**

---

## 4 DIFFUSION MODEL PARAMETERIZATION AND TRAINING LOSS

In this section, we discuss how to parameterize the denoising model  $\hat{\mathbf{x}}_\theta$ , and how to specify the reconstruction loss weight  $w(\lambda_t)$ . We assume a standard variance-preserving diffusion process for which  $\sigma_t^2 = 1 - \alpha_t^2$ , and we use a cosine schedule  $\alpha_t = \cos(0.5\pi t)$ , similar to that introduced by Nichol & Dhariwal (2021a).

Ho et al. (2020) and much of the following work choose to parameterize the denoising model through directly predicting  $\epsilon$  with a neural network  $\hat{\epsilon}_\theta(\mathbf{z}_t)$ , which implicitly models  $\mathbf{x}$  as  $\hat{\mathbf{x}}_\theta(\mathbf{z}_t) = \frac{1}{\alpha_t} (\mathbf{z}_t - \sigma_t \hat{\epsilon}_\theta(\mathbf{z}_t))$ . In this case, the training loss is also usually defined as mean squared error in the  $\epsilon$ -space:

$$L_\theta = \|\epsilon - \hat{\epsilon}_\theta(\mathbf{z}_t)\|_2^2 = \left\| \frac{1}{\sigma_t} (\mathbf{z}_t - \alpha_t \mathbf{x}) - \frac{1}{\sigma_t} (\mathbf{z}_t - \alpha_t \hat{\mathbf{x}}_\theta(\mathbf{z}_t)) \right\|_2^2 = \frac{\alpha_t^2}{\sigma_t^2} \|\mathbf{x} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t)\|_2^2, \quad (9)$$

which can thus equivalently be seen as a weighted reconstruction loss in  $\mathbf{x}$ -space, where the weighting function is given by  $w(\lambda_t) = \exp(\lambda_t)$ , for log signal-to-noise ratio  $\lambda_t = \log[\alpha_t^2 / \sigma_t^2]$ .

When we tried this standard specification of the denoising diffusion model, the progressive distillation procedure proposed in the last section did not work well. The reason is that this parameterization is not well suited for distillation: when training the original diffusion model, and at the start of progressive distillation, the model is evaluated at a wide range of signal-to-noise ratios  $\alpha_t^2/\sigma_t^2$ , but as distillation progresses we increasingly evaluate at lower and lower signal-to-noise ratios. As the signal-to-noise ratio goes to zero, the effect of small changes in the neural network output  $\hat{\epsilon}_\theta(\mathbf{z}_t)$  on the implied prediction in  $\mathbf{x}$ -space is increasingly amplified, since  $\hat{\mathbf{x}}_\theta(\mathbf{z}_t) = \frac{1}{\alpha_t}(\mathbf{z}_t - \sigma_t \hat{\epsilon}_\theta(\mathbf{z}_t))$  divides by  $\alpha_t$  which approaches zero. This is not much of a problem when taking many steps, since the effect of early missteps is limited by clipping of the  $\mathbf{z}_t$  iterates, and later updates can correct any mistakes, but it becomes increasingly important as we decrease the number of sampling steps. Eventually, if we distill all the way down to a single sampling step, the input to the model is only pure noise  $\epsilon$ , which corresponds to a signal-to-noise ratio of zero, i.e.  $\alpha_t = 0, \sigma_t = 1$ . At this extreme, the link between  $\epsilon$ -prediction and  $\mathbf{x}$ -prediction breaks down completely: observed data  $\mathbf{z}_t = \epsilon$  is no longer informative of  $\mathbf{x}$  and predictions  $\hat{\epsilon}_\theta(\mathbf{z}_t)$  no longer implicitly predict  $\mathbf{x}$ . Examining our reconstruction loss (equation 9), we see that the weighting function  $w(\lambda_t)$  gives zero weight to the reconstruction loss at this signal-to-noise ratio.

For distillation to work, we thus need to parameterize the diffusion model in a way for which the implied prediction  $\hat{\mathbf{x}}_\theta(\mathbf{z}_t)$  remains stable as  $\lambda_t = \log[\alpha_t^2/\sigma_t^2]$  varies. We tried the following options, and found all to work well with progressive distillation:

- Predicting  $\mathbf{x}$  directly.
- Predicting both  $\mathbf{x}$  and  $\epsilon$ , via separate output channels  $\{\tilde{\mathbf{x}}_\theta(\mathbf{z}_t), \tilde{\epsilon}_\theta(\mathbf{z}_t)\}$  of the neural network, and then merging the predictions via  $\hat{\mathbf{x}} = \sigma_t^2 \tilde{\mathbf{x}}_\theta(\mathbf{z}_t) + \alpha_t(\mathbf{z}_t - \sigma_t \tilde{\epsilon}_\theta(\mathbf{z}_t))$ , thus smoothly interpolating between predicting  $\mathbf{x}$  directly and predicting via  $\epsilon$ .
- Predicting  $\mathbf{v} \equiv \alpha_t \epsilon - \sigma_t \mathbf{x}$ , which gives  $\hat{\mathbf{x}} = \alpha_t \mathbf{z}_t - \sigma_t \hat{\mathbf{v}}_\theta(\mathbf{z}_t)$ , as we show in Appendix D.

In Section 5.1 we test all three parameterizations on training an original diffusion model (no distillation), and find them to work well there also.

In addition to determining an appropriate parameterization, we also need to decide a reconstruction loss weighting  $w(\lambda_t)$ . The setup of Ho et al. (2020) weights the reconstruction loss by the signal-to-noise ratio, implicitly gives a weight of zero to data with zero SNR, and is therefore not a suitable choice for distillation. We consider two alternative training loss weightings:

- $L_\theta = \max(\|\mathbf{x} - \hat{\mathbf{x}}_t\|_2^2, \|\epsilon - \hat{\epsilon}_t\|_2^2) = \max(\frac{\alpha_t^2}{\sigma_t^2}, 1)\|\mathbf{x} - \hat{\mathbf{x}}_t\|_2^2$ ; ‘truncated SNR’ weighting.
- $L_\theta = \|\mathbf{v}_t - \hat{\mathbf{v}}_t\|_2^2 = (1 + \frac{\alpha_t^2}{\sigma_t^2})\|\mathbf{x} - \hat{\mathbf{x}}_t\|_2^2$ ; ‘SNR+1’ weighting.

We examine both choices in our ablation study in Section 5.1, and find both to be good choices for training diffusion models. In practice, the choice of loss weighting also has to take into account how  $\alpha_t, \sigma_t$  are sampled during training, as this sampling distribution strongly determines the weight the expected loss gives to each signal-to-noise ratio. Our results are for a cosine schedule  $\alpha_t = \cos(0.5\pi t)$ , where time is sampled uniformly from  $[0, 1]$ . In Figure 2 we visualize the resulting loss weightings, both including and excluding the effect of the cosine schedule.

## 5 EXPERIMENTS

In this section we empirically validate the progressive distillation algorithm proposed in Section 3, as well as the parameterizations and loss weightings considered in Section 4. We consider various image generation benchmarks, with resolution varying from  $32 \times 32$  to  $128 \times 128$ . All experiments use the cosine schedule  $\alpha_t = \cos(0.5\pi t)$ , similar to that introduced by Nichol & Dhariwal (2021a), and all models use a U-Net architecture similar to that introduced by Ho et al. (2020), but with BigGAN-style up- and downsampling (Brock et al., 2019), as used in the diffusion modeling setting by Nichol & Dhariwal (2021a); Song et al. (2021c). Our training setup closely matches the open source code by Ho et al. (2020). Exact details are given in Appendix E. We plan to open source the code soon.

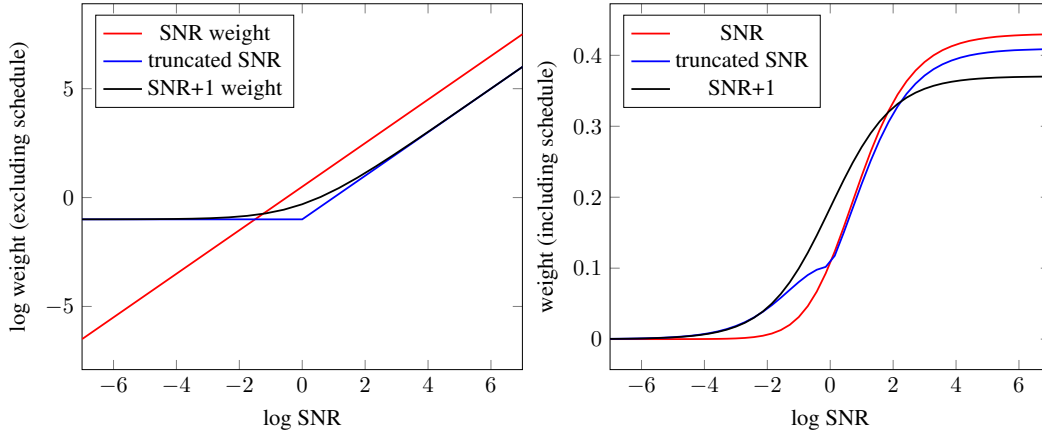


Figure 2: **Left:** Log weight assigned to reconstruction loss  $\|\mathbf{x} - \hat{\mathbf{x}}_\lambda\|_2^2$  as a function of the log-SNR  $\lambda = \log[\alpha^2/\sigma^2]$ , for each of our considered training loss weightings, excluding the influence of the  $\alpha_t, \sigma_t$  schedule. **Right:** Weights assigned to the reconstruction loss including the effect of the cosine schedule  $\alpha_t = \cos(0.5\pi t)$ , with  $t \sim U[0, 1]$ . The weights are only defined up to a constant, and we have adjusted these constants to fit this graph.

Network Output	Loss Weighting	Stochastic sampler	DDIM sampler
(x, $\epsilon$ ) combined	SNR	2.54/9.88	2.78/9.56
	Truncated SNR	2.47/9.85	2.76/9.49
	SNR+1	2.52/9.79	2.87/9.45
$\mathbf{x}$	SNR	2.65/9.80	2.75/9.56
	Truncated SNR	2.53/ <b>9.92</b>	<b>2.51/9.58</b>
	SNR+1	2.56/9.84	2.65/9.52
$\epsilon$	SNR	2.59/9.84	2.91/9.52
	Truncated SNR	N/A	N/A
	SNR+1	2.56/9.77	3.27/9.41
$\mathbf{v}$	SNR	2.65/9.86	3.05/9.56
	Truncated SNR	<b>2.45/9.80</b>	2.75/9.52
	SNR+1	2.49/9.77	2.87/9.43

Table 1: Generated sample quality as measured by FID and Inception Score (FID/IS) on unconditional CIFAR-10, training the original model (no distillation), and comparing different parameterizations and loss weightings discussed in Section 4. All reported results are averages over 3 random seeds of the best metrics obtained over 2 million training steps; nevertheless we find results are still  $\pm 0.1$  due to the noise inherent in training our models. Taking the neural network output to represent a prediction of  $\epsilon$  in combination with the Truncated SNR loss weighting leads to divergence.

### 5.1 MODEL PARAMETERIZATION AND TRAINING LOSS

As explained in Section 4, the standard method of having our model predict  $\epsilon$ , and minimizing mean squared error in the  $\epsilon$ -space (Ho et al., 2020), is not appropriate for use with progressive distillation. We therefore proposed various alternative parameterizations of the denoising diffusion model that are stable under the progressive distillation procedure, as well as various weighting functions for the reconstruction error in  $\mathbf{x}$ -space. Here, we perform a complete ablation experiment of all parameterizations and loss weightings considered in Section 4. For computational efficiency, and for comparisons to established methods in the literature, we use unconditional CIFAR-10 as the benchmark. We measure performance of undistilled models trained from scratch, to avoid introducing too many factors of variation into our analysis.

Table 1 lists the results of the ablation study. Overall results are fairly close across different parameterizations and loss weights, which is as expected given the similarity of the proposed losses as studied in Figure 2. All proposed stable model specifications achieve excellent performance, with the exception of the combination of outputting  $\epsilon$  with the neural network and weighting the loss with the truncated SNR, which we find to be unstable. If a single parameterization and loss weight is to be recommended, it is a good choice to predict  $x$  directly, in combination with the truncated SNR weight: this setting achieves the best FID and Inception scores when evaluated using the DDIM sampler.

## 5.2 PROGRESSIVE DISTILLATION

We evaluate our proposed progressive distillation algorithm on 4 data sets: CIFAR-10,  $64 \times 64$  downsampled ImageNet,  $128 \times 128$  LSUN bedrooms, and  $128 \times 128$  LSUN Church-Outdoor. For each data set we start by training a baseline model, after which we start the progressive distillation procedure. For CIFAR-10 we start progressive distillation from a teacher model taking 8192 steps. For the bigger data sets we start at 1024 steps. At every iteration of distillation we train for 50 thousand parameter updates, except for the distillation to 2 and 1 sampling steps, for which we use 100 thousand updates. For each model, the computational cost of progressive distillation to 4 sampling steps is comparable or less than for training the original model. We report FID results obtained after each iteration of the algorithm.

In Figure 4 we plot the resulting FID scores (Heusel et al., 2017) obtained for each number of sampling steps. We compare against the undistilled DDIM sampler, as well as to a highly optimized stochastic baseline sampler. For all four data sets progressive distillation produces near optimal results up to 4 or 8 sampling steps. At 2 or 1 sampling steps, the sample quality degrades relatively more quickly. In contrast, the quality of the DDIM and stochastic samplers degrades very sharply after reducing the number of sampling steps below 128. Overall, we conclude that progressive distillation is thus an attractive solution for computational budgets that allow less than or equal to 128 sampling steps.

Table 2 shows some of our results on CIFAR-10, and compares against other fast sampling methods in the literature: Our method compares favorably and attains higher sampling quality in fewer steps than most of the alternative methods. Figure 3 shows some random samples from our model obtained at different phases of the distillation process. Additional samples are provided in Appendix F



Figure 3: Random samples from our distilled  $64 \times 64$  ImageNet model, conditioned on the ‘malamute’ class, for varying numbers of sampling steps.

## 6 RELATED WORK ON FAST SAMPLING

Our proposed method is closest to the work of Luhman & Luhman (2021), who perform DDIM distillation in one step. A possible downside of their method is that it requires constructing a large data set by running the original model at its full number of sampling steps: the cost of distillation



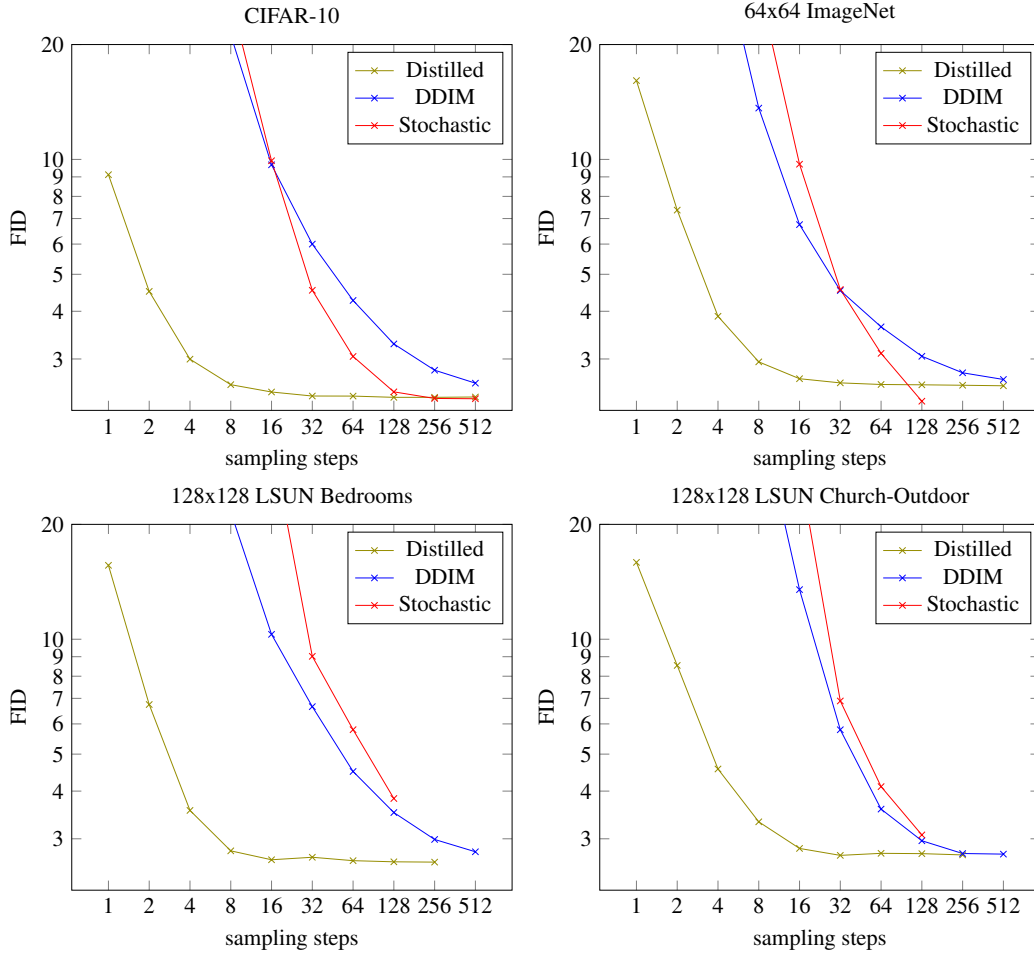


Figure 4: Sample quality results as measured by FID for our distilled model on unconditional CIFAR-10, class-conditional 64x64 ImageNet, 128x128 LSUN bedrooms, and 128x128 LSUN church-outdoor. We compare against the DDIM sampler and against an optimized stochastic sampler, each evaluated using the same models that were used to initialize the progressive distillation procedure. For CIFAR-10 we report an average over 4 random seeds. For the other data sets we only use a single run because of their computational demand. For the stochastic sampler we set the variance as a log-scale interpolation between an upper and lower bound on the variance, following Nichol & Dhariwal (2021b), but we use a single interpolation coefficient rather than a learned coefficient. We then tune this interpolation coefficient separately for each number of sampling steps and report only the best result for that number of steps: this way we obtained better results than with the learned interpolation.

thus scales linearly with this number of steps, which could potentially be prohibitive. In contrast, our method never needs to run the original model at the full number of sampling steps: at every distillation iteration, the number of model evaluations is independent of the number of teacher sampling steps, allowing our method to scale up to large numbers of teacher steps at a logarithmic cost in total distillation time.

DDIM (Song et al., 2021a) was originally shown to be effective for few-step sampling, as was the probability flow sampler (Song et al., 2021c). Jolicoeur-Martineau et al. (2021) study fast SDE integrators for reverse diffusion processes, and Tzen & Raginsky (2019b) study unbiased samplers which may be useful for fast, high quality sampling as well.

Other work on fast sampling can be viewed as manual or automated methods to adjust samplers or diffusion processes for fast generation. Nichol & Dhariwal (2021a); Kong & Ping (2021) describe



Method	Model evaluations	FID
Progressive Distillation (ours)	1	9.12
	2	4.51
	4	3.00
	8	2.57
Knowledge distillation (Luhman & Luhman, 2021)	1	9.36
DDIM (Song et al., 2021a)	10	13.36
	20	6.84
	50	4.67
	100	4.16
Dynamic step-size extrapolation + VP-deep (Jolicœur-Martineau et al., 2021)	48	82.42
	151	2.73
	180	2.44
	274	2.60
	330	2.56
FastDPM (Kong & Ping, 2021)	10	9.90
	20	5.05
	50	3.20
	100	2.86
Improved DDPM resampling (Nichol & Dhariwal, 2021a), our reimplementation	25	7.53
	50	4.99
LSGM (Vahdat et al., 2021)	138	2.10

Table 2: Comparison of fast sampling results on CIFAR-10 for diffusion models in the literature.

a methods to adjust discrete time diffusion model trained on many timesteps into models that can sample in few timesteps. Watson et al. (2021) describe a dynamic programming algorithm to reduce the number of timesteps for a diffusion model in a way that is optimal for log likelihood. Chen et al. (2021); Saharia et al. (2021); Ho et al. (2021) train diffusion models over continuous noise levels and tune samplers post training by adjusting the noise levels of a few-step discrete time reverse diffusion process. Their method is effective in highly conditioned settings such as text-to-speech and image super-resolution. San-Roman et al. (2021) train a new network to estimate the noise level of noisy data and show how to use this estimate to speed up sampling.

Alternative specifications of the diffusion model can also lend themselves to fast sampling, such as modified forward and reverse processes (Nachmani et al., 2021; Lam et al., 2021) and training diffusion models in latent space (Vahdat et al., 2021).

## 7 DISCUSSION

We have presented *progressive distillation*, a method to drastically reduce the number of sampling steps required for high quality generation of images, and potentially other data, using diffusion models with deterministic samplers like DDIM (Song et al., 2020). By making these models cheaper to run at test time, we hope to increase their usefulness for practical applications, for which running time and computational requirements often represent important constraints.

In the current work we limited ourselves to setups where the student model has the same architecture and number of parameters as the teacher model: in future work we hope to relax this constraint and explore settings where the student model is smaller, potentially enabling further gains in test time computational requirements. In addition, we hope to move past the generation of images and also explore progressive distillation of diffusion models for different data modalities such as e.g. audio (Chen et al., 2021).

In addition to the proposed distillation procedure, some of our progress was realized through different parameterizations of the diffusion model and its training loss. We expect to see more progress in this direction as the community further explores this model class.

## REPRODUCIBILITY STATEMENT

We provide full details on model architectures, training procedures, and hyperparameters in Appendix E, in addition to our discussion in Section 5. In Algorithm 2 we provide fairly detailed pseudocode that closely matches our actual implementation. We plan to open source our code and trained model checkpoints with the release of the final version of this paper.

## ETHICS STATEMENT

In general, generative models can have unethical uses, such as fake content generation, and they can suffer from bias if applied to data sets that are not carefully curated. The focus of this paper specifically is on speeding up generative models at test time in order to reduce their computational demands; we do not have specific concerns with regards to this contribution.

## REFERENCES

- Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. *CoRR*, abs/2107.03006, 2021.
- Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2019.
- Ruojin Cai, Guandao Yang, Hadar Averbuch-Elor, Zekun Hao, Serge Belongie, Noah Snively, and Bharath Hariharan. Learning gradient fields for shape generation. *arXiv preprint arXiv:2008.06520*, 2020.
- Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. WaveGrad: Estimating gradients for waveform generation. *International Conference on Learning Representations*, 2021.
- Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In *Advances in Neural Information Processing Systems*, pp. 6571–6583, 2018.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat GANs on image synthesis. *arXiv preprint arXiv:2105.05233*, 2021a.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. *CoRR*, abs/2105.05233, 2021b.
- Will Grathwohl, Ricky TQ Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. Ffjord: Free-form continuous dynamics for scalable reversible generative models. *arXiv preprint arXiv:1810.01367*, 2018.
- 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. In *Advances in Neural Information Processing Systems*, pp. 6626–6637, 2017.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, pp. 6840–6851, 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.
- Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Towards non-autoregressive language models. *arXiv preprint arXiv:2102.05379*, 2021.
- 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 P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *arXiv preprint arXiv:2107.00630*, 2021.
- Zhifeng Kong and Wei Ping. On fast sampling of diffusion probabilistic models. *arXiv preprint arXiv:2106.00132*, 2021.
- Max WY Lam, Jun Wang, Rongjie Huang, Dan Su, and Dong Yu. Bilateral denoising diffusion models. *arXiv preprint arXiv:2108.11514*, 2021.
- Haoying Li, Yifan Yang, Meng Chang, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *arXiv preprint arXiv:2104.14951*, 2021.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*, 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. *International Conference on Machine Learning*, 2021a.
- Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML*, 2021b.
- Chenhao Niu, Yang Song, Jiaming Song, Shengjia Zhao, Aditya Grover, and Stefano Ermon. Permutation invariant graph generation via score-based generative modeling. In *International Conference on Artificial Intelligence and Statistics*, pp. 4474–4484. PMLR, 2020.
- Prajit Ramachandran, Tom Le Paine, Pooya Khorrami, Mohammad Babaeizadeh, Shiyu Chang, Yang Zhang, Mark A Hasegawa-Johnson, Roy H Campbell, and Thomas S Huang. Fast generation for convolutional autoregressive models. *arXiv preprint arXiv:1704.06001*, 2017.
- Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *arXiv preprint arXiv:2104.07636*, 2021.
- 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, 2015.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *International Conference on Learning Representations*, 2021a.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, pp. 11895–11907, 2019.
- Yang Song and Stefano Ermon. Improved techniques for training score-based generative. *Advances in Neural Information Processing Systems*, 2020.
- Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *arXiv e-prints*, pp. arXiv–2101, 2021b.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*, 2021c.

- Yuxuan Song, Qiwei Ye, Minkai Xu, and Tie-Yan Liu. Discriminator contrastive divergence: Semi-amortized generative modeling by exploring energy of the discriminator. *arXiv preprint arXiv:2004.01704*, 2020.
- Belinda Tzen and Maxim Raginsky. Neural stochastic differential equations: Deep latent gaussian models in the diffusion limit. *arXiv preprint arXiv:1905.09883*, 2019a.
- Belinda Tzen and Maxim Raginsky. Theoretical guarantees for sampling and inference in generative models with latent diffusions. In *Conference on Learning Theory*, pp. 3084–3114. PMLR, 2019b.
- Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *arXiv preprint arXiv:2106.05931*, 2021.
- 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.