

# On Distillation of Guided Diffusion Models

Chenlin Meng<sup>1</sup>

chenlin@cs.stanford.edu

Robin Rombach<sup>2</sup>

robin@stability.ai

Ruiqi Gao<sup>3</sup>

ruiqiq@google.com

Diederik Kingma<sup>3</sup>

durk@google.com

Stefano Ermon<sup>1</sup>

ermon@cs.stanford.edu

Jonathan Ho<sup>3</sup>

jonathanho@google.com

Tim Salimans<sup>3</sup>

salimans@google.com

<sup>1</sup>Stanford University

<sup>2</sup>Stability AI & LMU Munich

<sup>3</sup>Google Research, Brain Team

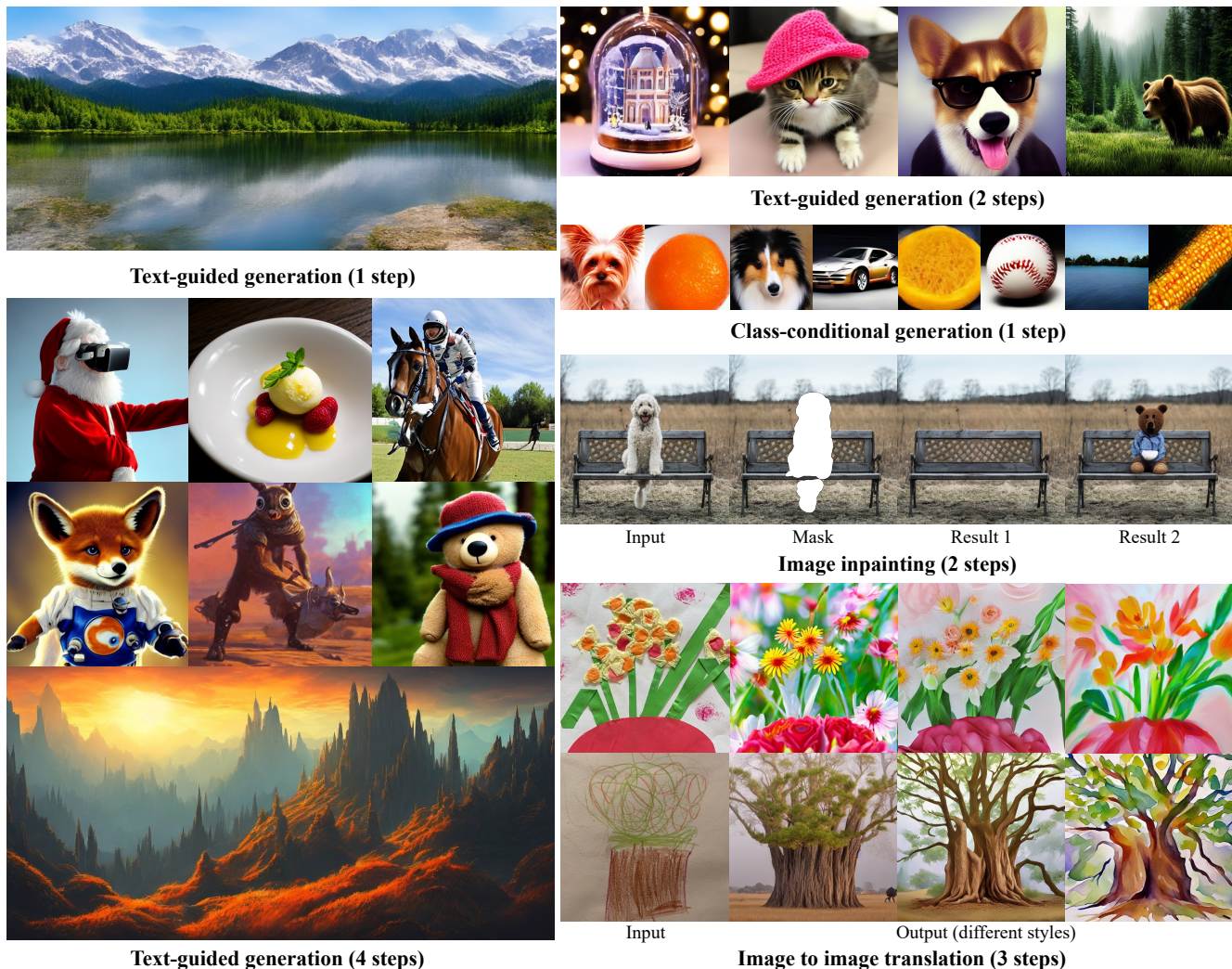


Figure 1. Distilled Stable Diffusion samples generated by our method. Our two-stage distillation approach is able to generate realistic images using only 1 to 4 denoising steps on various tasks. Compared to standard classifier-free guided diffusion models, we reduce the total number of sampling steps by at least 20 $\times$ .

## Abstract

Classifier-free guided diffusion models have recently been shown to be highly effective at high-resolution image generation, and they have been widely used in large-scale diffusion

frameworks including DALL-E 2, Stable Diffusion and Imagen. However, a downside of classifier-free guided diffusion models is that they are computationally expensive at inference time since they require evaluating two diffusion models, a class-conditional model and an unconditional model, tens to hundreds of times. To deal with this limitation, we pro-

\*Work partially done during an internship at Google

pose an approach to distilling classifier-free guided diffusion models into models that are fast to sample from: Given a pre-trained classifier-free guided model, we first learn a single model to match the output of the combined conditional and unconditional models, and then we progressively distill that model to a diffusion model that requires much fewer sampling steps. For standard diffusion models trained on the pixel-space, our approach is able to generate images visually comparable to that of the original model using as few as 4 sampling steps on ImageNet 64x64 and CIFAR-10, achieving FID/IS scores comparable to that of the original model while being up to 256 times faster to sample from. For diffusion models trained on the latent-space (e.g., Stable Diffusion), our approach is able to generate high-fidelity images using as few as 1 to 4 denoising steps, accelerating inference by at least 10-fold compared to existing methods on ImageNet 256x256 and LAION datasets. We further demonstrate the effectiveness of our approach on text-guided image editing and inpainting, where our distilled model is able to generate high-quality results using as few as 2-4 denoising steps.

## 1. Introduction

Denosing diffusion probabilistic models (DDPMs) [4, 37, 39, 40] have achieved state-of-the-art performance on image generation [22, 26–28, 31], audio synthesis [11], molecular generation [44], and likelihood estimation [10]. Classifier-free guidance [6] further improves the sample quality of diffusion models and has been widely used in large-scale diffusion model frameworks including GLIDE [23], Stable Diffusion [28], DALL-E 2 [26], and Imagen [31]. However, one key limitation of classifier-free guidance is its low sampling efficiency—it requires evaluating two diffusion models tens to hundreds of times to generate one sample. This limitation has hindered the application of classifier-free guidance models in real-world settings. Although distillation approaches have been proposed for diffusion models [33, 38], these approaches are not directly applicable to classifier-free guided diffusion models. To deal with this issue, we propose a two-stage distillation approach to improving the sampling efficiency of classifier-free guided models. In the first stage, we introduce a single student model to match the combined output of the two diffusion models of the teacher. In the second stage, we *progressively distill* the model learned from the first stage to a fewer-step model using the approach introduced in [33]. Using our approach, a *single* distilled model is able to handle a wide range of different guidance strengths, allowing for the trade-off between sample quality and diversity efficiently. To sample from our model, we consider existing deterministic samplers in the literature [33, 38] and further propose a stochastic sampling process.

Our distillation framework can not only be applied to standard diffusion models trained on the pixel-space [4, 36, 39],

but also diffusion models trained on the latent-space of an auto-encoder [28, 35] (e.g., Stable Diffusion [28]). For diffusion models directly trained on the pixel-space, our experiments on ImageNet 64x64 and CIFAR-10 show that the proposed distilled model can generate samples visually comparable to that of the teacher using only 4 steps and is able to achieve comparable FID/IS scores as the teacher model using as few as 4 to 16 steps on a wide range of guidance strengths (see Fig. 2). For diffusion model trained on the latent-space of an encoder [28, 35], our approach is able to achieve comparable visual quality to the base model using as few as 1 to 4 sampling steps (at least 10× fewer steps than the base model) on ImageNet 256×256 and LAION 512×512, matching the performance of the teacher (as evaluated by FID) with only 2-4 sampling steps. To the best of our knowledge, our work is the first to demonstrate the effectiveness of distillation for both pixel-space and latent-space classifier-free diffusion models. Finally, we apply our method to text-guided image inpainting and text-guided image editing tasks [20], where we reduce the total number of sampling steps to as few as 2-4 steps, demonstrating the potential of the proposed framework in style-transfer and image-editing applications [20, 41].



Figure 2. Class-conditional samples from our two-stage (deterministic) approach on ImageNet 64x64 for diffusion models trained on the pixel-space. By varying the guidance weight  $w$ , our distilled model is able to trade-off between sample diversity and quality, while achieving good results using as few as *one* sampling step.

## 2. Background on diffusion models

Given samples  $\mathbf{x}$  from a data distribution  $p_{\text{data}}(\mathbf{x})$ , noise scheduling functions  $\alpha_t$  and  $\sigma_t$ , we train a diffusion model  $\hat{\mathbf{x}}_{\theta}$ , with parameter  $\theta$ , via minimizing the weighted mean squared error [4, 36, 39, 40]

$$\mathbb{E}_{t \sim U[0,1], \mathbf{x} \sim p_{\text{data}}(\mathbf{x}), \mathbf{z}_t \sim q(\mathbf{z}_t | \mathbf{x})} [\omega(\lambda_t) \|\hat{\mathbf{x}}_{\theta}(\mathbf{z}_t) - \mathbf{x}\|_2^2], \quad (1)$$

where  $\lambda_t = \log[\alpha_t^2 / \sigma_t^2]$  is a signal-to-noise ratio [10],  $q(\mathbf{z}_t | \mathbf{x}) = \mathcal{N}(\mathbf{z}_t; \alpha_t \mathbf{x}, \sigma_t^2 \mathbf{I})$  and  $\omega(\lambda_t)$  is a pre-specified weighting function [10].

Once the diffusion model  $\hat{\mathbf{x}}_{\theta}$  is trained, one can use discrete-time DDIM sampler [38] to sample from the model. Specifically, the DDIM sampler starts with  $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$



and updates as follows

$$\mathbf{z}_s = \alpha_s \hat{\mathbf{x}}_\theta(\mathbf{z}_t) + \sigma_s \frac{\mathbf{z}_t - \alpha_t \hat{\mathbf{x}}_\theta(\mathbf{z}_t)}{\sigma_t}, \quad s = t - 1/N \quad (2)$$

with  $N$  the total number of sampling steps. The final sample will then be generated using  $\hat{\mathbf{x}}_\theta(\mathbf{z}_0)$ .

**Classifier-free guidance** Classifier-free guidance [6] is an effective approach shown to significantly improve the sample quality of class-conditioned diffusion models, and has been widely used in large-scale diffusion models including GLIDE [23], Stable Diffusion [28], DALL-E 2 [26] and Imagen [31]. Specifically, it introduces a guidance weight parameter  $w \in \mathbb{R}^{\geq 0}$  to trade-off between sample quality and diversity. To generate a sample, classifier-free guidance evaluates both a conditional diffusion model  $\hat{\mathbf{x}}_{c,\theta}$ , where  $c$  is the context (*e.g.*, class label, text prompt) to be conditioned on, and a jointly trained unconditional diffusion model  $\hat{\mathbf{x}}_\theta$  at each update step, using  $\hat{\mathbf{x}}_\theta^w = (1+w)\hat{\mathbf{x}}_{c,\theta} - w\hat{\mathbf{x}}_\theta$  as the model prediction in Eq. (2). As each sampling update requires evaluating two diffusion models, sampling with classifier-free guidance is often expensive [6].

**Progressive distillation** Our approach is inspired by *progressive distillation* [33], an effective method for improving the sampling speed of (unguided) diffusion models by repeated distillation. Until now, this method could not be directly applied to distilling classifier-free guided models or studied for samplers other than the deterministic DDIM sampler [33, 38]. In this paper we resolve these shortcomings.

**Latent diffusion models (LDMs)** [21, 24, 28, 35] increase the training and inference efficiency of diffusion models (directly learned on the pixel-space) by modeling images in the latent space of a pre-trained regularized autoencoder, where the latent representations are usually of lower dimensionality than the pixel-space. Latent diffusion models can be considered as an alternative to cascaded diffusion approaches [5], which rely on one or more super-resolution diffusion models to scale up a low-dimensional image to the desired target resolution.

In this work, we will apply our distillation framework to classifier-free guided diffusion models learned on both pixel-space [4, 36, 39] and latent-space [21, 24, 28, 35].

### 3. Distilling a guided diffusion model

In the following, we discuss our approach for distilling a classifier-free guided diffusion model [6] into a student model that requires fewer steps to sample from. Using a *single* distilled model conditioned on the guidance strength, our model can capture a wide range of classifier-free guidance levels, allowing for the trade-off between sample quality and diversity efficiently.

Given a trained guided model  $[\hat{\mathbf{x}}_{c,\theta}, \hat{\mathbf{x}}_\theta]$  (teacher) either on the pixel-space or latent-space, our approach can be decomposed into two stages.

#### 3.1. Stage-one distillation

In the first stage, we introduce a student model  $\hat{\mathbf{x}}_{\eta_1}(\mathbf{z}_t, w)$ , with learnable parameter  $\eta_1$ , to match the output of the teacher at any time-step  $t \in [0, 1]$ . The student model can either be a continuous-time model [40] or a discrete-time model [4, 38] depending on whether the teacher model is discrete or continuous. For simplicity, in the following discussion, we assume both the student and teacher models are continuous as the algorithm for discrete models is almost identical.

A key functionality of classifier-free guidance [6] is its ability to easily trade-off between sample quality and diversity, which is controlled by a “guidance strength” parameter. This property has demonstrated utility in real-world applications [6, 23, 26, 28, 31], where the optimal “guidance strength” is often a user preference. Thus, we would also want our distilled model to maintain this property. Given a range of guidance strengths  $[w_{\min}, w_{\max}]$  we are interested in, we optimize the student model using the following objective

$$\mathbb{E}_{w \sim p_w, t \sim U[0,1], \mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[ \omega(\lambda_t) \|\hat{\mathbf{x}}_{\eta_1}(\mathbf{z}_t, w) - \hat{\mathbf{x}}_\theta^w(\mathbf{z}_t)\|_2^2 \right], \quad (3)$$

where  $\hat{\mathbf{x}}_\theta^w(\mathbf{z}_t) = (1+w)\hat{\mathbf{x}}_{c,\theta}(\mathbf{z}_t) - w\hat{\mathbf{x}}_\theta(\mathbf{z}_t)$ ,  $\mathbf{z}_t \sim q(\mathbf{z}_t|\mathbf{x})$  and  $p_w(w) = U[w_{\min}, w_{\max}]$ . Note that here, our distilled model  $\hat{\mathbf{x}}_{\eta_1}(\mathbf{z}_t, w)$  is also conditioned on the context  $c$  (*e.g.*, text prompt), but we drop the notation  $c$  in the paper for simplicity. We provide the detailed training algorithm in Algorithm 1 in the supplement.

To incorporate the guidance weight  $w$ , we introduce a  $w$ -conditioned model, where  $w$  is fed as an input to the student model. To better capture the feature, we apply Fourier embedding to  $w$ , which is then incorporated into the diffusion model backbone in a way similar to how the time-step was incorporated in [10, 33]. As initialization plays a key role in the performance [33], we initialize the student model with the same parameters as the conditional model of the teacher, except for the newly introduced parameters related to  $w$ -conditioning. The model architecture we use is a U-Net model similar to the ones used in [6] for pixel-space diffusion models and [1, 28] for latent-space diffusion models. We use the same number of channels and attention as used in [6] and the open-sourced Stable Diffusion repository\* for our experiments. We provide more details in the supplement.

#### 3.2. Stage-two distillation

In the second stage, we consider a discrete time-step scenario and progressively distill the learned model from the first-stage  $\hat{\mathbf{x}}_{\eta_1}(\mathbf{z}_t, w)$  into a fewer-step student model  $\hat{\mathbf{x}}_{\eta_2}(\mathbf{z}_t, w)$  with learnable parameter  $\eta_2$ , by halving the number of sampling steps each time. Letting  $N$  denote the number of sampling steps, given  $w \sim U[w_{\min}, w_{\max}]$  and

\*<https://github.com/CompVis/stable-diffusion>

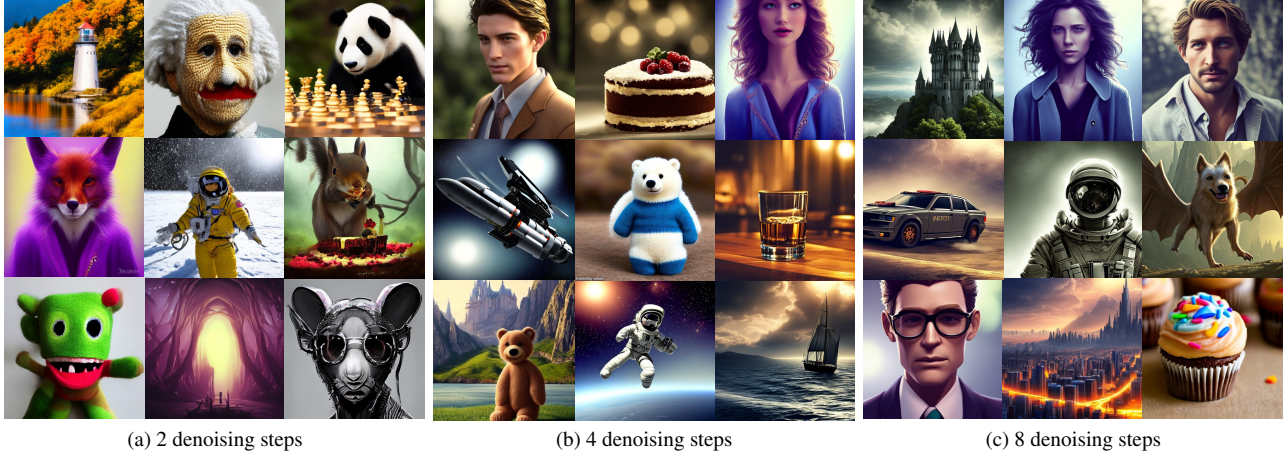


Figure 4. Text-guided generation on LAION (512x512) using our distilled Stable Diffusion model. Our model is able to generate high-quality image samples using 2, 4 or 8 denoising steps, significantly improving the inference efficiency of Stable Diffusion.

Method	$w = 0$		$w = 0.3$		$w = 1$		$w = 4$	
	FID ( $\downarrow$ )	IS ( $\uparrow$ )	FID ( $\downarrow$ )	IS ( $\uparrow$ )	FID ( $\downarrow$ )	IS ( $\uparrow$ )	FID ( $\downarrow$ )	IS ( $\uparrow$ )
Ours 1-step (D/S)	22.74 / 26.91	25.51 / 23.55	14.85 / 18.48	37.09 / 33.30	7.54 / 8.92	75.19 / 67.80	18.72 / <b>17.85</b>	157.46 / 148.97
Ours 4-step (D/S)	4.14 / 3.91	46.64 / 48.92	2.17 / 2.24	69.64 / 73.73	7.95 / 8.51	128.98 / 135.36	26.45 / 27.33	207.45 / 216.56
Ours 8-step (D/S)	2.79 / 2.44	50.72 / 55.03	<b>2.05</b> / 2.31	76.01 / 83.00	9.33 / 10.56	136.47 / 147.39	26.62 / 27.84	203.47 / <b>219.89</b>
Ours 16-step (D/S)	2.44 / <b>2.10</b>	52.53 / <b>57.81</b>	2.20 / 2.56	79.47 / <b>87.50</b>	9.99 / 11.63	139.11 / <b>153.17</b>	26.53 / 27.69	204.13 / 218.70
Single- $w$ 1-step	19.61	24.00	11.70	36.95	<b>6.64</b>	74.41	19.857	170.69
Single- $w$ 4-step	4.79	38.77	2.34	62.08	8.23	118.52	27.75	219.64
Single- $w$ 8-step	3.39	42.13	2.32	68.76	9.69	125.20	27.67	218.08
Single- $w$ 16-step	2.97	43.63	2.56	70.97	10.34	127.70	27.40	216.52
DDIM 16x2-step [38]	7.68	37.60	5.33	60.83	9.53	112.75	21.56	195.17
DDIM 32x2-step [38]	5.03	40.93	7.47	9.33	9.26	126.22	23.03	213.23
DDIM 64x2-step [38]	3.74	43.16	5.52	9.51	9.53	133.17	23.64	217.88
Teacher (DDIM 1024x2-step)	2.92	44.81	2.36	74.83	9.84	139.50	23.94	224.74

Table 1. ImageNet 64x64 distillation results for pixel-space diffusion models ( $w = 0$  refers to non-guided models). For our method,  $D$  and  $S$  stand for deterministic and stochastic sampler respectively. We observe that training the model conditioned on a guidance interval  $w \in [0, 4]$  performs comparably with training a model on a fixed  $w$  (see Single- $w$ ). Our approach significantly outperforms DDIM when using fewer steps, and is able to match the teacher performance using as few as 8 to 16 steps.

$t \in \{1, \dots, N\}$ , we train the student model to match the output of two-step DDIM sampling of the teacher (i.e., from  $t/N$  to  $t - 0.5/N$  and from  $t - 0.5/N$  to  $t - 1/N$ ) in one step, following the approach of [33]. After distilling the  $2N$  steps in the teacher model to  $N$  steps in the student model, we can use the  $N$ -step student model as the new teacher model, repeat the same procedure, and distill the teacher model into a  $N/2$ -step student model. At each step, we initialize the student model with the parameters of the teacher. We provide the training algorithm and extra details in the supplementary material.

### 3.3. $N$ -step deterministic and stochastic sampling

Once the model  $\hat{\mathbf{x}}_{\eta_2}$  is trained, given a specified guidance strength  $w \in [w_{\min}, w_{\max}]$ , we can perform sampling via the DDIM update rule in Eq. (2). We note that given the distilled model  $\hat{\mathbf{x}}_{\eta_2}$ , this sampling procedure is *deterministic* given the initialization  $\mathbf{z}_1^w$ . In fact, we can also perform  $N$ -step *stochastic* sampling: We apply one deterministic sampling step with two-times the original step-length (i.e.,

the same as a  $N/2$ -step deterministic sampler) and then perform one stochastic step backward (i.e., perturb with noise) using the original step-length, a process inspired by [9]. With  $\mathbf{z}_1^w \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , we use the following update rule when  $t > 1/N$

$$\mathbf{z}_k^w = \alpha_k \hat{\mathbf{x}}_{\eta_2}(\mathbf{z}_t^w) + \sigma_k \frac{\mathbf{z}_t^w - \alpha_t \hat{\mathbf{x}}_{\eta_2}(\mathbf{z}_t^w)}{\sigma_t}, \quad (4)$$

$$\text{where } \mathbf{z}_s^w = (\alpha_s / \alpha_k) \mathbf{z}_k^w + \sigma_{s|k} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); \quad (5)$$

$$\mathbf{z}_h^w = \alpha_h \hat{\mathbf{x}}_{\eta_2}(\mathbf{z}_s^w) + \sigma_h \frac{\mathbf{z}_s^w - \alpha_s \hat{\mathbf{x}}_{\eta_2}(\mathbf{z}_s^w)}{\sigma_s}, \quad (6)$$

$$\text{where } \mathbf{z}_k^w = (\alpha_k / \alpha_h) \mathbf{z}_h^w + \sigma_{k|h} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (7)$$

In the above equations,  $h = t - 3/N$ ,  $k = t - 2/N$ ,  $s = t - 1/N$  and  $\sigma_{a|b}^2 = (1 - e^{\lambda_a - \lambda_b}) \sigma_a^2$ . When  $t = 1/N$ , we use deterministic update Eq. (2) to obtain  $\mathbf{z}_0^w$  from  $\mathbf{z}_{1/N}^w$ . We provide an illustration of the process in Fig. 5, where the number of denoising steps is 4. We note that compared to the *deterministic* sampler, performing *stochastic* sampling requires evaluating the model at slightly different time-steps,



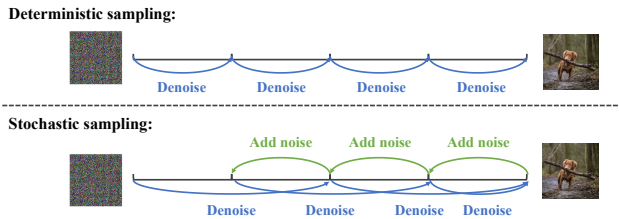


Figure 5. Sampling procedures of the distilled model where the number of denoising steps is 4.

and would require small modifications to training algorithm for the edge cases. We provide the algorithm and more details in the supplementary material.

## 4. Experiments

In this section, we evaluate the performance of our distillation approach on pixel-space diffusion models (*i.e.* DDPM [4]) and latent-space diffusion models (*i.e.* Stable Diffusion [28]). We further apply our approach to text-guided image editing and inpainting tasks. Experiments show that our approach is able to achieve competitive performance while using as few as 2-4 steps on all tasks.

### 4.1. Distillation for pixel-space guided models

In this experiment, we consider class-conditional diffusion models directly trained on the pixel-space [4, 6, 33].

**Settings** We focus on ImageNet 64x64 [30] and CIFAR-10 [12] as higher-resolution image generation in this scenario often relies on combining with other super-resolution techniques [5, 31]. We explore different ranges for the guidance weight and observe that all ranges work comparably and therefore use  $[w_{min}, w_{max}] = [0, 4]$  for the experiments. The baselines we consider include DDPM ancestral sampling [4] and DDIM [38]. The teacher model we use is a 1024x2-step DDIM model, where the conditional and unconditional components both use 1024 DDIM denoising steps. To better understand how the guidance weight  $w$  should be incorporated, we also include models trained using a single fixed  $w$  as a baseline. We use the same pre-trained teacher model for all the methods for fair comparisons. Following [4, 6, 39], we use a U-Net [29, 39] architecture for the baselines, and the same U-Net backbone with the introduced  $w$ -embedding for our two-step student models (see Sec. 3). Following [33], we use a  $v$ -prediction model for both datasets.

**Results** We report the performance as evaluated in FID [3] and Inception scores (IS) [32] for all approaches on ImageNet 64x64 in Fig. 6 and Tab. 1 and provide extended ImageNet 64x64 and CIFAR-10 results in the supplement. We observe that our distilled model is able to match a teacher guided DDIM model with 1024x2 sampling steps using only 4-16 steps, achieving a speedup for up to 256 $\times$ . We emphasize that, using our approach, a *single* distilled model is able to match the teacher performance on a wide range

of guidance strengths. This has not been achieved by any previous methods.

### 4.2. Distillation for latent-space guided models

After demonstrating the effectiveness of our method on pixel-space class-guided diffusion models in Sec. 4.1, we now expand its scope to latent-space diffusion models. In the following sections, we show the effectiveness of our approach on Latent Diffusion [28] on a variety of tasks, including class-conditional generation, text-to-image generation, image inpainting and text-guided style-transfer [20].

In the following experiments, we use the open-sourced latent-space diffusion models [28] as the teacher models. As  $v$ -prediction teacher model tends to perform better than  $\epsilon$ -prediction model, we *fine-tune* the open-sourced  $\epsilon$ -prediction models into  $v$ -prediction teacher models. We provide more details in the supplementary material.

#### 4.2.1 Class-conditional generation

In this section, we apply our method to a class-conditional latent diffusion model pre-trained on ImageNet  $256 \times 256$ . We start from the DDIM teacher model with 512 sampling steps, and use the output as the target to train our distilled model. We use a batch size of 512 and uniformly sample the guidance strength  $w \in [w_{min} = 0, w_{max} = 14]$  during training.

**Results** Empirically, we find that our distilled model is able to match the performance of the teacher model (originally trained on 1000 steps) in terms of FID scores while using only 2 or 4 sampling steps. We also achieve significantly better performance than DDIM when using 1-4 sampling steps (see Fig. 11). Qualitatively, we find that samples synthesized using a single denoising step still yield satisfying results, while the baseline fails to generate images with meaningful contents. We provide extra samples in the supplementary material.

Similar to the pixel-based results in Fig. 6, we also observe the trade-off between sampling quality and diversity as measured by FID and Inception Score for our distilled latent diffusion model. Following Kynkäänniemi *et al* [13], we further compute improved precision and recall metrics for this experiment in the appendix.

#### 4.2.2 Text-guided image generation

In this section, we focus on the text-guided Stable Diffusion model pretrained on subsets<sup>†</sup> of LAION-5B [34] at a resolution  $512 \times 512$ . We then follow our two-stage approach introduced in Sec. 3 and distill the guided model in 3000 gradient updates into a  $w$ -conditioned model using

<sup>†</sup>[https://github.com/CompVis/stable-diffusion/blob/main/Stable\\_Diffusion\\_v1\\_Model\\_Card.md](https://github.com/CompVis/stable-diffusion/blob/main/Stable_Diffusion_v1_Model_Card.md)

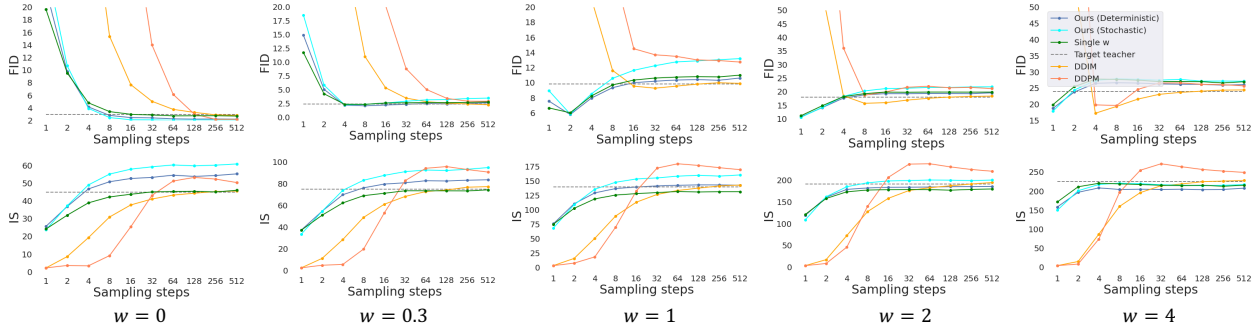


Figure 6. ImageNet 64x64 sample quality evaluated by FID and IS scores. Our distilled model significantly outperform the DDPM and DDIM baselines, and is able to match the performance of the teacher using as few as 8 to 16 steps. By varying  $w$ , a *single* distilled model is able to capture the trade-off between sample diversity and quality.



Figure 7. Text-guided Stable Diffusion results. We distill the public *Stable Diffusion* model using the proposed pipeline, arriving at a model that achieves high sample quality using only four denoising steps (*left*). When sampling from the original model using four DDIM steps, the generated samples have clear artifacts (*middle*). When using eight DDIM steps, the results get better (*right*), but are still blurry and less consistent than the distilled results using fewer steps. More samples are provided in Fig. 4.

$w \in [w_{min} = 2, w_{max} = 14]$ , and a batch size of 512. Although we can condition on a broader range of  $w$  for the distilled (student) model, the utility remains unclear as we typically do not exceed the normal guidance range when sampling with the teacher model. The final model is obtained by applying progressive distillation for 2000 training steps per stage, except when for the low-step regime of 1, 2, and 4 steps, where we train for 20000 gradient updates. A detailed analysis of the convergence properties of this model in the supplement.

Method	2-step	4-step	8-step
DPM [16]	98.9/0.20	34.3/0.29	31.7/0.32
DPM++ [18]	98.8/0.20	34.1/0.29	25.6/0.32
Ours	37.3/0.27	26.0/0.30	26.9/0.30

Table 2. FID/CLIP scores on LAION 512X512 ( $w = 8.0$ ). We point out that DPM and DPM++ use both the conditional and unconditional components for sampling. Depending on the implementation, this either requires higher peak memory or two times more sampling steps.

**Results** We present samples in Fig. 4. We evaluate the resulting model both qualitatively and quantitatively. For the latter analysis, we follow [31] and evaluate CLIP [25] and FID scores to assess text-image alignment and quality, respectively. We use the open-sourced ViT-g/14 [7] CLIP model for evaluation. The quantitative results in Fig. 10 show that our method can significantly increase the performance in both metrics over DDIM sampling on the base model for

2 and 4 sampling steps. For 8 steps, these metrics do not show a significant difference. However, when considering the corresponding samples in Fig. 7 we can observe a stark difference in terms of visual image quality. In contrast to the 8-step DDIM samples from the original model, the distilled samples are sharper and more coherent. We hypothesize that FID and CLIP do not fully capture these differences in our evaluation setting on COCO2017 [14], where we used 5000 random captions from the validation set. We further compute the FID and CLIP scores for our distilled LAION 512x512 model and compare them with the DPM [16] and DPM++ [18] solver in Tab. 2. We observe that our method is able to achieve significantly better performance when the denoising step is 2 or 4. Furthermore, we stress that stage-one of our method already decreases the number of function evaluations by a factor of 2, as we distill the classifier-free guidance step into a single model. Depending on the exact implementation (batched vs. sequential network evaluation), this either decreases peak memory or sampling time compared to existing solvers [16, 18, 38].

### 4.2.3 Text-guided image-to-image translation

In this section, we perform experiments on text-guided image-to-image translation with SDEdit [20] using our distilled model from Sec. 4.2.2. Following SDEdit [20], we perform stochastic encoding in the latent space, but instead



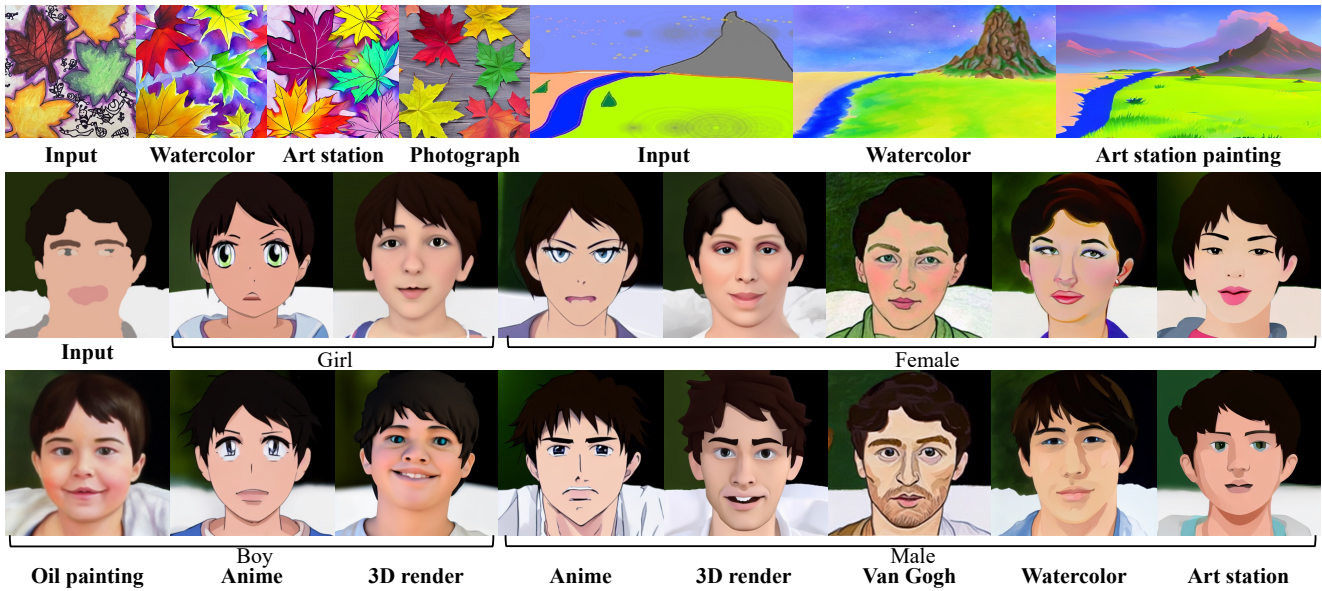


Figure 8. Text-guided image-to-image translation [20] with the distilled Stable Diffusion model (3 denoising steps). We observe that our model is able to generate high-quality and faithful outputs using only 3 denoising steps.



Figure 9. Image inpainting with our distilled Stable Diffusion model (4 denoising steps). Our model is able to generate high-quality image inpainting results using 4 denoising steps on unseen data.

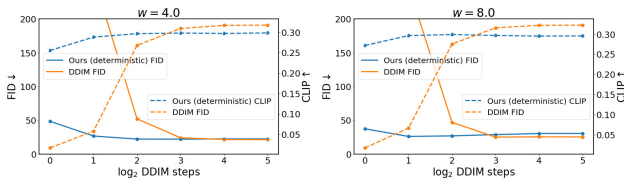


Figure 10. FID and CLIP ViT-g/14 score for text-to-image generation at  $512 \times 512$  px using the distilled Stable Diffusion model. The results are evaluated on 5000 captions from the COCO2017 [14] validation set. Our distilled latent diffusion model is able to generate high-quality image samples using significantly less sampling steps than the original model while achieving similar or better FID and CLIP scores, especially in the low-step regime.

use the deterministic sampler of the distilled model to perform deterministic decoding. We consider input image and text of various kinds and provide qualitative results in Fig. 8. We observe that our distilled model generates high-quality style-transfer results using as few as 3 denoising steps. We provide more analysis on the trade-off between sample quality, controllability and efficiency in the supplement.

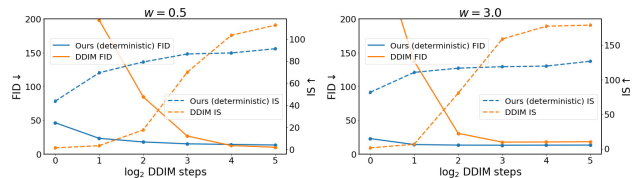


Figure 11. FID and Inception Score for class-conditional image generation on ImageNet ( $256 \times 256$ ) with distilled latent diffusion. The results are evaluated on 5000 samples. Our distilled latent diffusion model is able to generate high-quality image samples using significantly less sampling steps (up to a factor of 16) than the original model while achieving similar or better FID scores.

#### 4.2.4 Image inpainting

In this section, we apply our approach to a pre-trained image inpainting latent diffusion model. We use the open-source *Stable Diffusion Inpainting*<sup>‡</sup> image-inpainting model. This model is a fine-tuned version of the pure text-to-image *Stable Diffusion* model from above, where additional input channels

<sup>‡</sup><https://huggingface.co/runwayml/stable-diffusion-inpainting>



Figure 12. Style transfer comparison on ImageNet 64x64 for pixel-space models. For our approach, we use a distilled encoder and decoder. For the baseline, we encode and decode using DDIM. We use  $w = 0$  and 16 sampling steps for both the encoder and decoder. We observe that our method achieves more realistic outputs.

were added to process masks and masked images.

We use the same distillation algorithm as used in the previous section. For training, we start from the  $v$ -prediction teacher model sampled with 512 DDIM steps, and use the output as the target to optimize our student model. We present qualitative results in Fig. 9, demonstrating the potential of our method for fast, real-world image editing applications. For additional training details and a quantitative evaluation, see the supplementary.

### 4.3. Progressive distillation for encoding

In this experiment, we explore distilling the encoding process for the teacher model and perform experiments on style-transfer in a setting similar to [41]. We focus on pixel-space diffusion models pre-trained on ImageNet  $64 \times 64$ . Specifically, to perform style-transfer between two domains  $A$  and  $B$ , we encode the image from domain- $A$  using a diffusion model trained on domain- $A$ , and then decode with a diffusion model trained on domain- $B$ . As the encoding process can be understood as reversing the DDIM sampling process, we perform distillation for both the encoder and decoder with classifier-free guidance, and compare with a DDIM encoder and decoder in Fig. 12. We also explore how modifying the guidance strength  $w$  can impact the performance and provide more details in the supplementary material.

## 5. Related Work

Our approach is related to existing works on improving the sampling speed of diffusion models [4, 37, 40]. For instance, denoising diffusion implicit models (DDIM [38]), probability flow sampler [40], fast SDE integrators [8] have been proposed to improve the sampling speed of diffusion models. Other works develop higher-order solvers [17], exponential integrators [15], and dynamic programming based approach [43] to accelerating sampling speed. However, none of these approaches have achieved comparable performance as our method on distilling classifier-free guided diffusion models.

Existing distillation-based methods for diffusion models are mainly designed for non-classifier-free guided diffusion

models. For instance, [19] proposes to predict the data from noise in one single step by inverting a deterministic encoding of DDIM, [2] proposes to achieve faster sampling speed by distilling higher order solvers into an additional prediction head of the neural network backbone [2]. *Progressive distillation* [33] is perhaps the most relevant work. Specifically, it proposes to progressively distill a pre-trained diffusion model into a fewer-step student model with the same model architecture. However, none of these approaches are directly applicable or have been applied to classifier-free guided diffusion models. They are also unable to capture a range of different guidance strengths using one single distilled model. On the contrary, by incorporating the guidance strength into the model architecture and training the model using a two-stage procedure, our approach is able to match the performance of the teacher model on a wide range of guidance strength using one *single* model. Using our method, one single model can capture the trade-off between sample quality and diversity, enabling the real-world application of classifier-free guided diffusion models, where the guidance strength is often specified by users. Moreover, none of the above distillation approaches have been applied to or shown effectiveness for latent-space text-to-image models. Finally, most fast sampling approaches [33, 38, 40] only consider using deterministic sampling schemes to improve the sampling speed. In this work, we further develop an effective stochastic sampling approach to sample from the distilled models.

## 6. Conclusion

In this paper, we propose a distillation approach for guided diffusion models [6]. Our two-stage approach allows us to significantly speed up popular but relatively inefficient guided diffusion models. We show that our approach can reduce the inference cost of classifier-free guided pixel-space and latent-space diffusion models by at least an order of magnitude. Empirically, we show that our approach is able to produce visually appealing results with only 2 steps, achieving a comparable FID score to the teacher with as few as 4 to 8 steps. We further demonstrate practical applications of our distillation approach to text-guided image-to-image translation and inpainting tasks. We hope that by significantly reducing the inference cost of classifier-free guided diffusion models, our method will promote creative applications as well as the wider adoption of image generation systems. In the future work, we aim to further improve the performance in the two and one sampling step regimes.

## Acknowledgements

We thank the anonymous reviewers for their insightful discussions and feedback. All experiments on Stable Diffusion are supported by Stability AI.



## References

- [1] Prafulla Dhariwal and Alex Nichol. Diffusion models beat GANs on image synthesis. *arXiv preprint arXiv:2105.05233*, 2021. 3
- [2] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Genie: Higher-order denoising diffusion solvers. In *Advances in Neural Information Processing Systems*, 2022. 8
- [3] 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*, pages 6626–6637, 2017. 5
- [4] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, pages 6840–6851, 2020. 2, 3, 5, 8
- [5] 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. 3, 5
- [6] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 2, 3, 5, 8, 11
- [7] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. If you use this software, please cite it as below. 6, 24
- [8] 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. 8
- [9] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *arXiv preprint arXiv:2206.00364*, 2022. 4
- [10] Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *arXiv preprint arXiv:2107.00630*, 2021. 2, 3, 11, 22, 23, 25
- [11] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. DiffWave: A Versatile Diffusion Model for Audio Synthesis. *International Conference on Learning Representations*, 2021. 2
- [12] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5
- [13] Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *Advances in Neural Information Processing Systems*, 32, 2019. 5, 24
- [14] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 6, 7, 24
- [15] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. In *International Conference on Learning Representations*, 2022. 8
- [16] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps, 2022. 6, 23, 24
- [17] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *arXiv:2206.00927*, 2022. 8, 11
- [18] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models, 2022. 6, 11, 13, 14, 23, 24, 28
- [19] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021. 8
- [20] Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 2, 5, 6, 7, 25, 26, 27
- [21] Gautam Mittal, Jesse Engel, Curtis Hawthorne, and Ian Simon. Symbolic music generation with diffusion models. *arXiv preprint arXiv:2103.16091*, 2021. 3
- [22] Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *International Conference on Machine Learning*, 2021. 2
- [23] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 2, 3
- [24] Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn. Diffusion autoencoders: Toward a meaningful and decodable representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10619–10629, 2022. 3
- [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 6, 24
- [26] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 2, 3
- [27] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 2
- [28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2, 3, 5
- [29] 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*, pages 234–241. Springer, 2015. 5
- [30] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 5
- [31] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 2, 3, 5, 6
- [32] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. In *Advances in Neural Information Processing Systems*, pages 2234–2242, 2016. 5
- [33] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 2, 3, 4, 5, 8, 11, 19, 22, 23, 25
- [34] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022. 5
- [35] Abhishek Sinha, Jiaming Song, Chenlin Meng, and Stefano Ermon. D2c: Diffusion-decoding models for few-shot conditional generation. *Advances in Neural Information Processing Systems*, 34:12533–12548, 2021. 2, 3
- [36] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265, 2015. 2, 3
- [37] Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *arXiv preprint arXiv:1503.03585*, March 2015. 2, 8
- [38] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *International Conference on Learning Representations*, 2021. 2, 3, 4, 5, 6, 8, 11, 13, 14, 16, 22, 24, 28
- [39] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, pages 11895–11907, 2019. 2, 3, 5
- [40] 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*, 2021. 2, 3, 8
- [41] Xuan Su, Jiaming Song, Chenlin Meng, and Stefano Ermon. Dual diffusion implicit bridges for image-to-image translation. *arXiv preprint arXiv:2203.08382*, 2022. 2, 8, 12
- [42] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2149–2159, 2022. 25
- [43] Daniel Watson, William Chan, Jonathan Ho, and Mohammad Norouzi. Learning fast samplers for diffusion models by differentiating through sample quality. In *International Conference on Learning Representations*, 2022. 8
- [44] Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geometric diffusion model for molecular conformation generation. *arXiv preprint arXiv:2203.02923*, 2022. 2
- [45] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 24



## A. Results overview

In this section, we provide an overview table for the speed-up we achieved for pixel-space and latent-space diffusion models (see Tab. 3). We also provide extra samples from the text-guided image generation model as well as comparison with DDIM [38], DPM [17] and DPM++ [18] solvers in Fig. 13 and Fig. 14. We provide more experimental details on pixel-space distillation in Appendix B and latent-space distillation in Appendix C.

## B. Pixel-space distillation

### B.1. Teacher model

The model architecture we use is a U-Net model similar to the ones used in [6]. The model is parameterized to predict  $\mathbf{v}$  as discussed in [33]. We use the same training setting as [6].

### B.2. Stage-one distillation

The model architecture we use is a U-Net model similar to the ones used in [6]. We use the same number of channels and attention as used in [6] for both ImageNet 64x64 and CIFAR-10. As mentioned in Section 3, we also make the model take  $w$  as input. Specifically, we apply Fourier embedding to  $w$  before combining with the model backbone. The way we incorporate  $w$  is the same as how time-step is incorporated to the model as used in [10, 33]. We parameterize the model to predict  $\mathbf{v}$  as discussed in [33]. We train the distilled model using Algorithm 1. We train the model using SNR loss [10, 33]. For ImageNet 64x64, we use learning rate  $3e-4$ , with EMA decay 0.9999; for CIFAR-10, we use learning rate  $1e-3$ , with EMA decay 0.9999. We initialize the student model with parameters from the teacher model except for the parameters related to  $w$ -embedding.

---

#### Algorithm 1 Stage-one distillation

---

**Require:** Trained classifier-free guidance teacher model  $[\hat{\mathbf{x}}_{c,\theta}, \hat{\mathbf{x}}_\theta]$

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

**Require:** Loss weight function  $\omega(\cdot)$

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

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

$t \sim U[0, 1]$  ▷ Sample time

$w \sim U[w_{\min}, w_{\max}]$  ▷ Sample guidance

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

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

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

$\hat{\mathbf{x}}_\theta^w(\mathbf{z}_t) = (1 + w)\hat{\mathbf{x}}_{c,\theta}(\mathbf{z}_t) - w\hat{\mathbf{x}}_\theta(\mathbf{z}_t)$  ▷ Compute target

$L_{\eta_1} = \omega(\lambda_t) \|\hat{\mathbf{x}}_\theta^w(\mathbf{z}_t) - \hat{\mathbf{x}}_{\eta_1}(\mathbf{z}_t, w)\|_2^2$  ▷ Loss

$\eta_1 \leftarrow \eta_1 - \gamma \nabla_{\eta_1} L_{\eta_1}$  ▷ Optimization

**end while**

---

### B.3. Stage-two distillation for deterministic sampler

We use the same model architectures as the ones used in Stage-one (see Appendix B.2). We train the distilled model using Algorithm 2. We first use the student model from Stage-one as the teacher model. We start from 1024 DDIM sampling steps and progressively distill the student model from Stage-one to a one step model. We train the student model for 50,000 parameter updates, except for sampling step equals to one or two where we train the model for 100,000 parameter updates, before the number of sampling step is halved and the student model becomes the new teacher model. At each sampling step, we initialize the student model with the parameters from the teacher model. We train the model using SNR truncation loss [10, 33]. For each step, we linearly anneal the learning rate from  $1e-4$  to 0 during each parameter update. We do not use EMA decay for training. Our training setting follows the setting in [33] closely.

---

#### Algorithm 2 Stage-two distillation for deterministic sampler

---

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

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

**Require:** Loss weight function  $\omega(\cdot)$

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

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

$\eta_2 \leftarrow \eta$  ▷ Init student from teacher

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

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

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

$w \sim U[w_{\min}, w_{\max}]$  ▷ Sample guidance

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

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

# 2 steps of DDIM with teacher

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

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

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

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

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

$L_{\eta_2} = \omega(\lambda_t) \|\tilde{\mathbf{x}}^w - \hat{\mathbf{x}}_{\eta_2}(\mathbf{z}_t, w)\|_2^2$

$\eta_2 \leftarrow \eta_2 - \gamma \nabla_{\eta_2} L_{\eta_2}$

**end while**

$\eta \leftarrow \eta_2$  ▷ Student becomes next teacher

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

**end for**

---

### B.4. Stage-two distillation for stochastic sampling

We train the distilled model using Algorithm 3. We use the same model architecture and training setting as Stage-two distillation described in Appendix B.3 for both ImageNet 64x64 and CIFAR-10: The main difference here is that our distillation target corresponds to taking a sampling step that

Space	Task	Dataset	Metric	Student diffusion step	Comparable teacher diffusion step	Speed-up
Pixel-space	class-conditional generation	CIFAR-10	FID	4	1024 DDIM×2	×512
	class-conditional generation	CIFAR-10	IS	4	1024 DDIM×2	×512
	class-conditional generation	ImageNet 64×64	FID	8	1024 DDIM×2	×256
	class-conditional generation	ImageNet 64×64	IS	8	1024 DDIM×2	×256
Latent-space	class-conditional generation	ImageNet 256×256	FID	2	16 DDIM ×2	×16
	class-conditional generation	ImageNet 256×256	Recall	2	16 DDIM ×2	×16
	text-guided generation	LAION-5B 512× 512	FID	2	16 DDIM / 8 DPM++ ×2	×16 / ×8
	text-guided generation	LAION-5B 512× 512	CLIP	4	8 DDIM / 4 DPM++ ×2	×8 / ×4

Table 3. Speed-up overview for pixel-space diffusion and latent-space diffusion. We note that the original model (without distillation) requires evaluating both the unconditional and the conditional diffusion model at each denoising step. Our model, on the other hand, only requires evaluating one diffusion model at each denoising step. This is because in our stage-one distillation, we distill the output of the unconditional and conditional models into the output of one model. Thus our method further decreases either the peak memory or sampling time by a half compared to the original model.

is twice as large as for the deterministic sampler. We provide visualization for samples with varying guidance strengths  $w$  in Fig. 15.

---

**Algorithm 3** Stage-two distillation for stochastic sampler

---

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

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

**Require:** Loss weight function  $\omega(\cdot)$

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

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

$\eta_2 \leftarrow \eta$  ▷ Init student from teacher

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

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

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

$w \sim U[w_{\min}, w_{\max}]$  ▷ Sample guidance

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

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

**if**  $t > 1/N$  **then**

# 2 steps of DDIM with teacher

$t' = t - 1/N, t'' = t - 2/N$

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

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

$\alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}^w, w))$

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

**else** ▷ Edge case

# 1 step of DDIM with teacher

$t' = t - 1/N$

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

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

**end if**

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

$L_{\eta_2} = \omega(\lambda_t) \|\tilde{\mathbf{x}}^w - \hat{\mathbf{x}}_{\eta_2}(\mathbf{z}_t, w)\|_2^2$

$\eta_2 \leftarrow \eta_2 - \gamma \nabla_{\eta_2} L_{\eta_2}$

**end while**

$\eta \leftarrow \eta_2$  ▷ Student becomes next teacher

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

**end for**

---

## B.5. Baseline samples

We provide extra samples for the DDIM baseline in Fig. 16 and Fig. 17.

## B.6. Extra distillation results

We provide the FID and IS results for our method and the baselines on ImageNet 64x64 and CIFAR-10 in Fig. 22b, Fig. 22a and Tab. 4. We also visualize the FID and IS trade-off curves for both datasets in Fig. 18 and Fig. 19, where we select guidance strength  $w = \{0, 0.3, 1, 2, 4\}$  for ImageNet 64x64 and  $w = \{0, 0.1, 0.2, 0.3, 0.5, 0.7, 1, 2, 4\}$  for CIFAR-10.

## B.7. Style transfer

We focus on ImageNet 64x64 for this experiment. As discussed in [41], one can perform style-transfer between domain A and B by encoding (performing reverse DDIM) an image using a diffusion model train on domain A and then decoding using DDIM with a diffusion model trained on domain B. We train the model using Algorithm 4. We use the same  $w$ -conditioned model architecture and training setting as discussed in Appendix B.3.



Figure 13. Text-guided image generation on LAION-5B ( $512 \times 512$ ). We compare our distilled model with the original model sampled with DDIM [38] and DPM++ [18]. We observe that our model, when using only two steps, is able to generate more realistic and higher quality images compared to the baselines using more steps. We note that both DDIM and DPM-Solver require evaluating both a conditional and an unconditional diffusion model at each denoising step, while we distill the two models into one model at our stage-one distillation and only require evaluating one model at each denoising step. Depending on the implementation, DDIM and DPM-Solver require either extra  $\times 2$  peak memory or  $\times 2$  sampling steps compared to our approach.



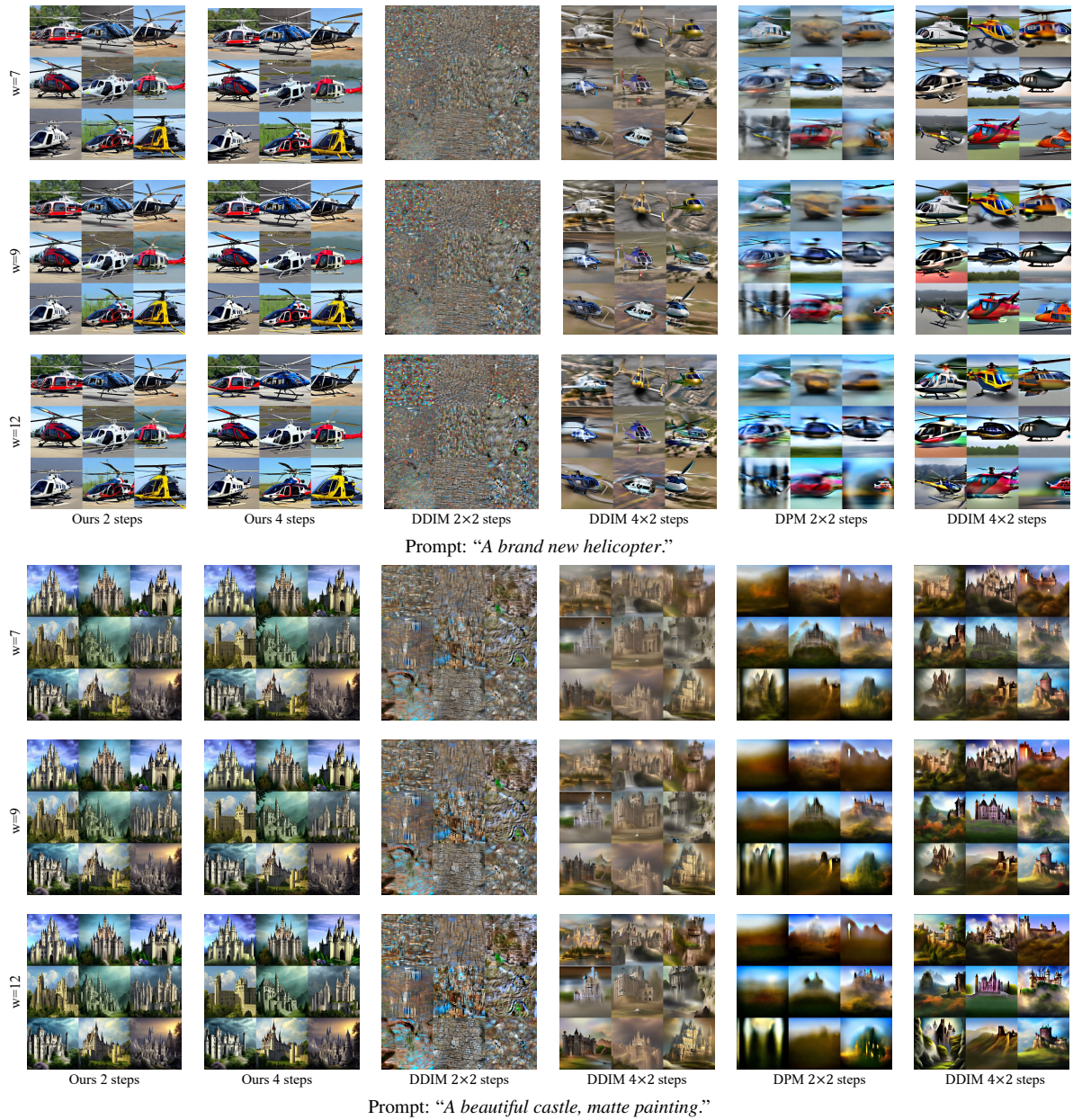


Figure 14. Text-guided image generation on LAION-5B ( $512 \times 512$ ). We compare our distilled model with the original model sampled with DDIM [38] and DPM++ [18]. We observe that our model, when using only two steps, is able to generate more realistic and higher quality images compared to the baselines using more steps. We note that both DDIM and DPM-Solver require evaluating both a conditional and an unconditional diffusion model at each denoising step, while we distill the two models into one model at our stage-one distillation and only require evaluating one model at each denoising step. Depending on the implementation, DDIM and DPM-Solver require either extra  $\times 2$  peak memory or  $\times 2$  sampling steps compared to our approach.



Figure 15. Class-conditional samples from our two-step (stochastic) approach on ImageNet 64x64. By varying the guidance weight  $w$ , our distilled model is able to trade-off between sample diversity and quality, while achieving visually pleasant results using as few as *one* sampling step.



Figure 16. ImageNet 64x64 class-conditional generation using DDIM (baseline)  $8 \times 2$  sampling steps. We observe clear artifacts when  $w = 0$ .



Figure 17. ImageNet 64x64 class-conditional generation using DDIM (baseline)  $16 \times 2$  sampling steps.



		ImageNet 64x64		CIFAR-10	
Guidance $w$	Model	FID ( $\downarrow$ )	IS ( $\uparrow$ )	FID ( $\downarrow$ )	IS ( $\uparrow$ )
$w = 0.0$	Ours 1-step (D/S)	22.74 / 26.91	25.51 / 23.55	8.34 / 10.65	8.63 / 8.42
	Ours 2-step (D/S)	9.75 / 10.67	36.69 / 37.12	4.48 / 4.81	9.23 / 9.30
	Ours 4-step (D/S)	4.14 / 3.91	46.64 / 48.92	3.18 / 3.28	9.50 / 9.60
	Ours 8-step (D/S)	2.79 / 2.44	50.72 / 55.03	2.86 / 3.11	9.68 / 9.74
	Ours 16-step (D/S)	2.44 / 2.10	52.53 / 57.81	2.78 / 3.12	9.67 / 9.76
	Single- $w$ 1-step	19.61	24.00	6.64	8.88
	Single- $w$ 4-step	4.79	38.77	3.14	9.47
	Single- $w$ 8-step	3.39	42.13	2.86	9.67
	Single- $w$ 16-step	2.97	43.63	2.75	9.65
	DDIM 16 $\times$ 2-step [38]	7.68	37.60	10.11	8.81
	DDIM 32 $\times$ 2-step [38]	5.03	40.93	6.67	9.17
	DDIM 64 $\times$ 2-step [38]	3.74	43.16	4.64	9.32
	Target (DDIM 1024 $\times$ 2-step)	2.92	44.81	2.73	9.66
	$w = 0.3$	Ours 1-step (D/S)	14.85 / 18.48	37.09 / 33.30	7.34 / 9.38
Ours 2-step (D/S)		5.052 / 5.81	54.44 / 54.37	4.23 / 4.74	9.45 / 9.45
Ours 4-step (D/S)		2.17 / 2.24	69.64 / 73.73	3.58 / 3.95	9.73 / 9.77
Ours 8-step (D/S)		2.05 / 2.31	76.01 / 83.00	3.54 / 3.96	9.87 / 9.90
Ours 16-step (D/S)		2.20 / 2.56	79.47 / 87.50	3.57 / 4.17	9.89 / 9.97
Single- $w$ 1-step		11.70	36.95	5.98	9.13
Single- $w$ 4-step		2.34	62.08	3.58	9.75
Single- $w$ 8-step		2.32	68.76	3.57	9.85
Single- $w$ 16-step		2.56	70.97	3.61	9.88
DDIM 16 $\times$ 2-step		5.33	60.83	10.83	8.96
DDIM 32 $\times$ 2-step		3.45	68.03	7.47	9.33
DDIM 64 $\times$ 2-step		2.80	72.55	5.52	9.51
Target (DDIM 1024 $\times$ 2-step)		2.36	74.83	3.65	9.83
$w = 1.0$		Ours 1-step (D/S)	7.54 / 8.92	75.19 / 67.80	8.62 / 10.27
	Ours 2-step (D/S)	5.77 / 5.83	109.97 / 108.38	6.88 / 7.52	9.64 / 9.55
	Ours 4-step (D/S)	7.95 / 8.51	128.98 / 135.36	7.39 / 7.64	9.86 / 9.87
	Ours 8-step (D/S)	9.33 / 10.56	136.47 / 147.39	7.81 / 7.85	9.9 / 10.05
	Ours 16-step (D/S)	9.99 / 11.63	139.11 / 153.17	7.97 / 8.34	10.00 / 10.05
	Single- $w$ 1-step	6.64	74.41	8.18	9.32
	Single- $w$ 4-step	8.23	118.52	7.66	9.88
	Single- $w$ 8-step	9.69	125.20	8.09	9.89
	Single- $w$ 16-step	10.34	127.70	8.30	9.95
	DDIM 16 $\times$ 2-step	9.53	112.75	14.81	8.98
	DDIM 32 $\times$ 2-step	9.26	126.22	11.44	9.36
	DDIM 64 $\times$ 2-step	9.53	133.17	9.79	9.64
	Target (DDIM 1024 $\times$ 2-step)	9.84	139.50	7.80	9.96
	$w = 2.0$	Ours 1-step (D/S)	10.71 / 10.55	118.55 / 108.37	13.23 / 14.33
Ours 2-step (D/S)		14.08 / 14.18	160.04 / 161.43	12.58 / 12.57	9.51 / 9.48
Ours 4-step (D/S)		17.61 / 18.23	178.29 / 184.45	13.83 / 13.24	9.70 / 9.77
Ours 8-step (D/S)		18.80 / 20.25	181.53 / 193.49	14.41 / 13.67	9.77 / 9.87
Ours 16-step (D/S)		19.25 / 21.11	183.17 / 197.71	14.80 / 14.28	9.79 / 9.84
Single- $w$ 1-step		11.12	120.74	13.31	9.23
Single- $w$ 4-step		18.14	172.74	14.04	9.70
Single- $w$ 8-step		19.24	176.74	14.67	9.77
Single- $w$ 16-step		19.81	177.69	15.04	9.79
DDIM 16 $\times$ 2-step		15.92	157.67	20.25	8.97
DDIM 32 $\times$ 2-step		16.85	175.72	17.27	9.29
DDIM 64 $\times$ 2-step		17.53	182.11	15.66	9.48
Target (DDIM 1024-step)		17.97	190.56	13.60	9.81
$w = 4.0$		Ours 1-step (D/S)	18.72 / 17.85	157.46 / 148.97	23.20 / 23.79
	Ours 2-step (D/S)	23.74 / 24.34	196.05 / 200.11	23.41 / 22.75	9.16 / 9.11
	Ours 4-step (D/S)	26.45 / 27.33	207.45 / 216.56	25.11 / 23.62	9.23 / 9.33
	Ours 8-step (D/S)	26.62 / 27.84	203.47 / 219.89	25.94 / 23.98	9.26 / 9.55
	Ours 16-step (D/S)	26.53 / 27.69	204.13 / 218.70	26.01 / 24.40	9.33 / 9.50
	Single- $w$ 1-step	19.857	170.69	23.17	8.93
	Single- $w$ 4-step	27.75	219.64	24.45	9.32
	Single- $w$ 8-step	27.67	218.08	24.83	9.38
	Single- $w$ 16-step	27.40	216.52	25.11	9.37
	DDIM 16 $\times$ 2-step	21.56	195.17	27.99	8.71
	DDIM 32 $\times$ 2-step	23.03	213.23	25.07	9.07
	DDIM 64 $\times$ 2-step	23.64	217.88	23.41	9.17
	Target (DDIM 1024 $\times$ 2-step)	23.94	224.74	21.28	9.54

Table 4. Distillation results on ImageNet 64x64 and CIFAR-10 ( $w = 0$  refers to non-guided models). For our method,  $D$  and  $S$  stand for deterministic and stochastic sampler respectively. We observe that training the model conditioned on an guidance interval  $w \in [0, 4]$  performs comparably with training a model on a fixed  $w$  (see Single- $w$ ). Our approach significantly outperforms DDIM when using fewer steps, and is able to match the teacher performance using as few as 8 to 16 steps. We also note that DDIM and DDPM evaluates both an unconditional and a conditional diffusion model at each denoising step, giving rise to the  $\times 2$  overhead either for peak memory or sampling steps.



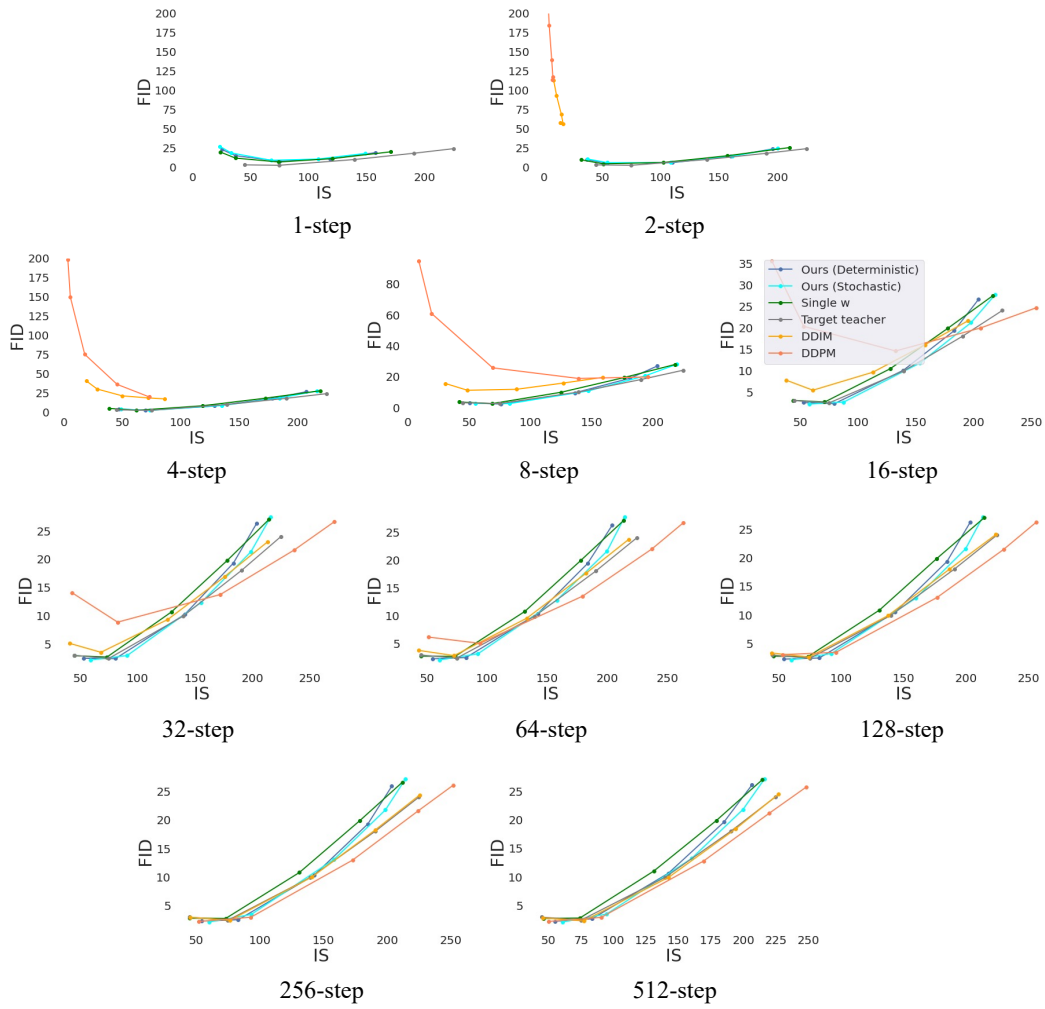


Figure 18. FID and IS score trade-off on ImageNet 64x64. We plot the results using guidance strength  $w = \{0, 0.3, 1, 2, 4\}$ . For the 1-step plot, the curves of DDIM and DDPM are too far to be visualized.

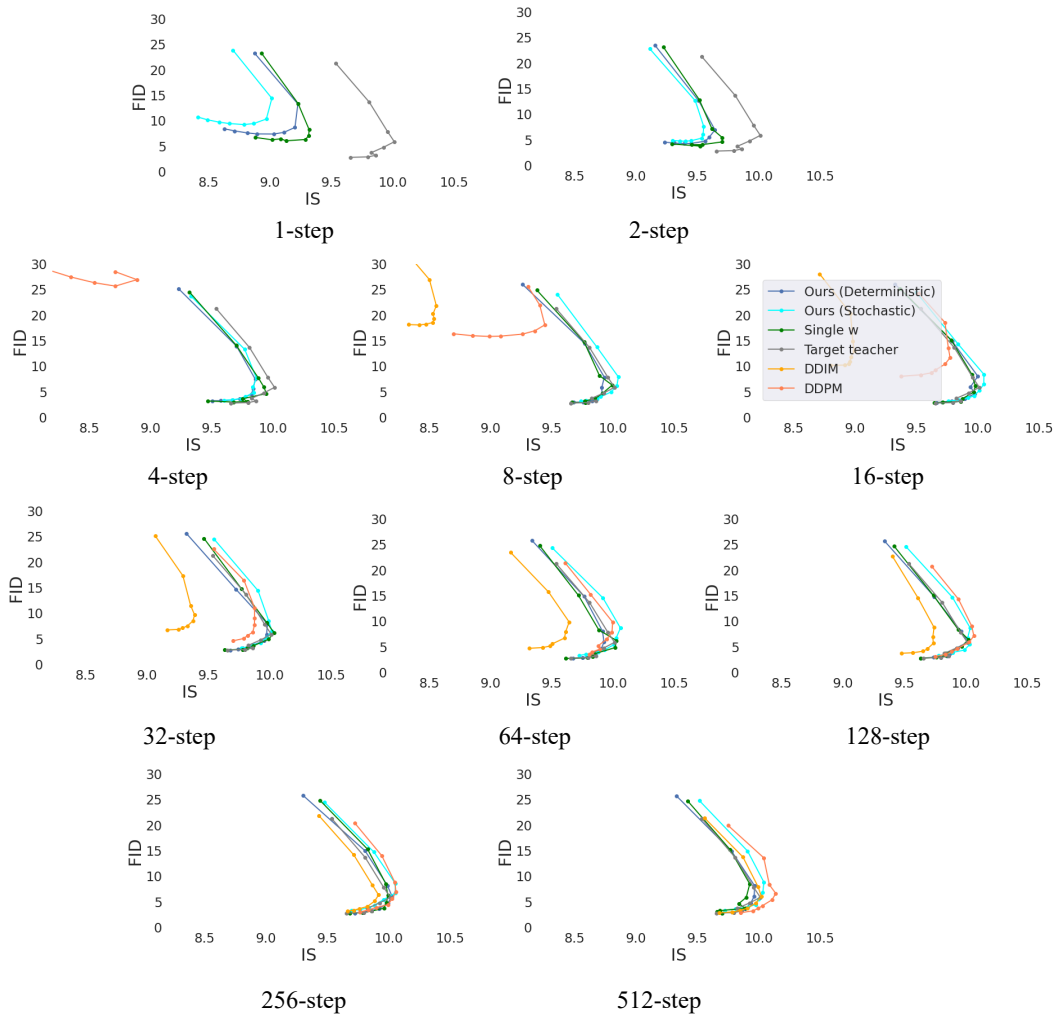


Figure 19. FID and IS score trade-off on CIFAR-10. We plot the results using guidance strength  $w = \{0, 0.1, 0.2, 0.3, 0.5, 0.7, 1, 2, 4\}$ . For the 1-step and 2-step plots, the curves of DDIM and DDPM are too far away to be visualized. For the 4-step plot, the curve of DDIM is too far away to be visualized.



Figure 20. Style transfer on ImageNet 64x64 for pixel-space models (orange to bell pepper). We use a distilled 16-step encoder and decoder. We fix the encoder guidance strength to be 0 and vary the decoder guidance strength from 0 to 4. As we increase  $w$ , we notice a trade-off between sample diversity and sharpness.

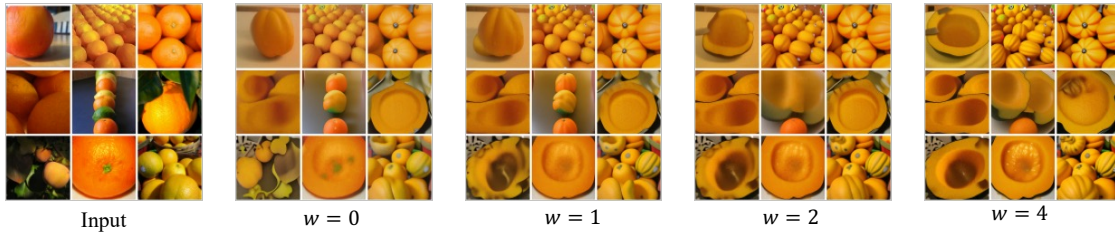


Figure 21. Style transfer on ImageNet 64x64 (orange to acorn squash). We use a distilled 16-step encoder and decoder. We fix the encoder guidance strength to be 0 and vary the decoder guidance strength from 0 to 4. As we increase the guidance strength  $w$ , we notice a trade-off between sample diversity and sharpness.

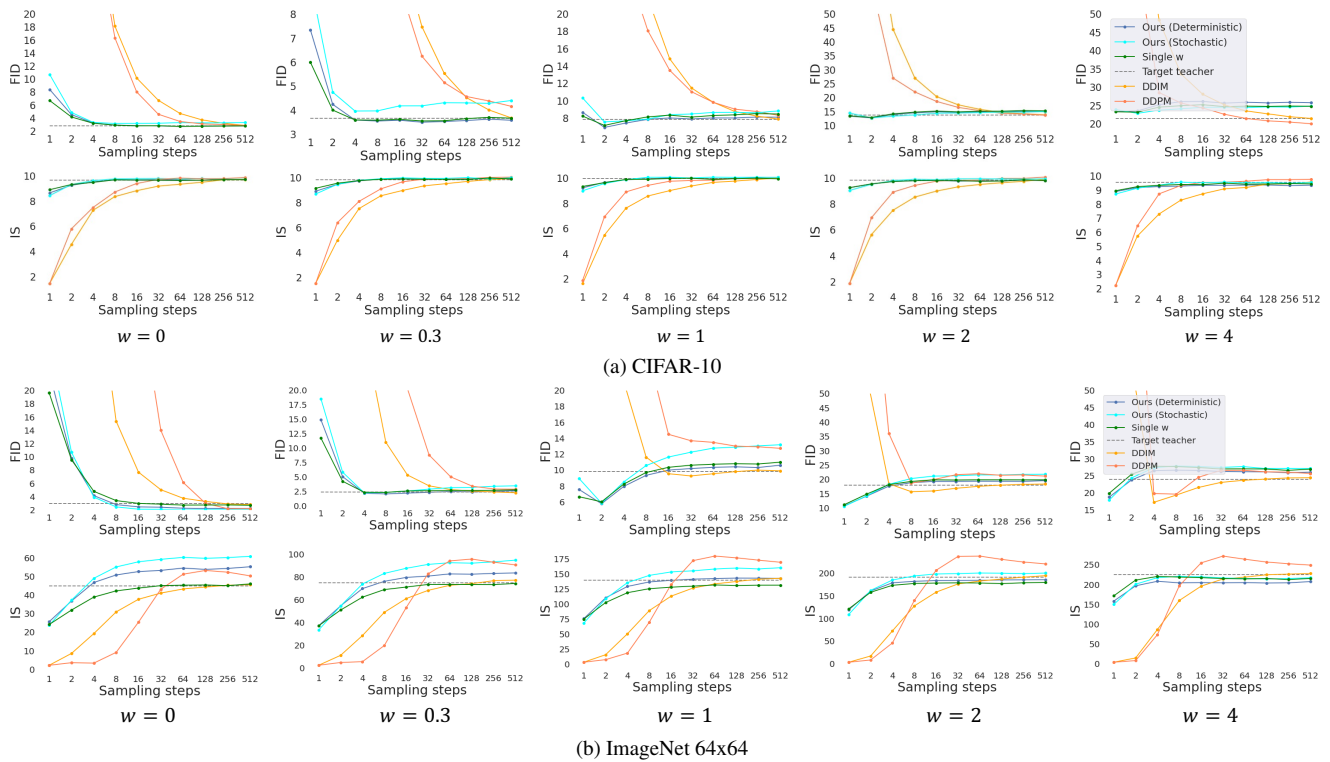


Figure 22. CIFAR-10 and ImageNet sample quality evaluated by FID and IS scores for pixel-space diffusion models. We follow the setting of [33] for our evaluation. We note that, the DDPM and DDIM baseline require evaluating both an unconditional and a conditional diffusion model at each denoising step for classifier-free guidance, giving rise to either an extra  $\times 2$  overhead for peak memory or an extra  $\times 2$  sampling steps than the “Sampling steps” value shown in the plot. Our distilled model significantly outperform the DDPM and DDIM baselines, and is able to match the performance of the teacher using as few as 4 to 16 steps. By varying  $w$ , a *single* distilled model is able to capture the trade-off between sample diversity and quality.



---

**Algorithm 4** Encoder distillation

---

**Require:** Trained teacher model  $\hat{\mathbf{x}}_\eta(\mathbf{z}_t, w)$ **Require:** Data set  $\mathcal{D}$ **Require:** Loss weight function  $\omega(\cdot)$ **Require:** Student sampling steps  $N$ **for**  $K$  iterations **do** $\eta_2 \leftarrow \eta$   $\triangleright$  Init student from teacher**while** not converged **do** $\mathbf{x} \sim \mathcal{D}$  $t = i/N, i \sim \text{Cat}[0, 1, \dots, N-1]$  $w \sim U[w_{\min}, w_{\max}]$   $\triangleright$  Sample guidance $\epsilon \sim N(0, I)$  $\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$ 

# 2 steps of reversed DDIM with

teacher

 $t' = t + 0.5/N, t'' = t + 1/N$  $\mathbf{z}_{t'}^w = \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_t, w) + \frac{\sigma_{t'}}{\sigma_t} (\mathbf{z}_t - \alpha_t \hat{\mathbf{x}}_\eta(\mathbf{z}_t, w))$  $\mathbf{z}_{t''}^w = \alpha_{t''} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}^w, w) + \frac{\sigma_{t''}}{\sigma_{t'}} (\mathbf{z}_{t'}^w - \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}^w, w))$  $\tilde{\mathbf{x}}^w = \frac{\mathbf{z}_{t''}^w - (\sigma_{t''}/\sigma_t) \mathbf{z}_t}{\alpha_{t''} - (\sigma_{t''}/\sigma_t) \alpha_t}$   $\triangleright$  Teacher  $\hat{\mathbf{x}}$  target $\lambda_t = \log[\alpha_t^2/\sigma_t^2]$  $L_{\eta_2} = \omega(\lambda_t) \|\tilde{\mathbf{x}}^w - \hat{\mathbf{x}}_{\eta_2}(\mathbf{z}_t, w)\|_2^2$  $\eta_2 \leftarrow \eta_2 - \gamma \nabla_{\eta_2} L_{\eta_2}$ **end while** $\eta \leftarrow \eta_2$   $\triangleright$  Student becomes next teacher $N \leftarrow N/2$   $\triangleright$  Halve number of sampling steps**end for**

---

---

**Algorithm 5** Two-student progressive distillation

---

**Require:** Trained classifier-free guidance teacher model  $[\hat{\mathbf{x}}_{c,\theta}, \hat{\mathbf{x}}_\theta]$ **Require:** Data set  $\mathcal{D}$ **Require:** Loss weight function  $\omega(\cdot)$ **Require:** Student sampling steps  $N$ **for**  $K$  iterations **do** $\eta \leftarrow \theta$   $\triangleright$  Init student from teacher**while** not converged **do** $\mathbf{x} \sim \mathcal{D}$  $t = i/N, i \sim \text{Cat}[1, 2, \dots, N]$  $w \sim U[w_{\min}, w_{\max}]$   $\triangleright$  Sample guidance $\epsilon \sim N(0, I)$  $\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$  $\hat{\mathbf{x}}_\theta^w(\mathbf{z}_t) = (1+w) \hat{\mathbf{x}}_{c,\theta}(\mathbf{z}_t) - w \hat{\mathbf{x}}_\theta(\mathbf{z}_t)$   $\triangleright$ 

Compute target

# 2 steps of DDIM with teacher

 $t' = t - 0.5/N, t'' = t - 1/N$  $\mathbf{z}_{t'}^w = \alpha_{t'} \hat{\mathbf{x}}_\theta^w(\mathbf{z}_t) + \frac{\sigma_{t'}}{\sigma_t} (\mathbf{z}_t - \alpha_t \hat{\mathbf{x}}_\theta^w(\mathbf{z}_t))$  $\mathbf{z}_{c,t''}^w = \alpha_{t''} \hat{\mathbf{x}}_{c,\theta}(\mathbf{z}_{t'}^w) + \frac{\sigma_{t''}}{\sigma_{t'}} (\mathbf{z}_{t'}^w - \alpha_{t'} \hat{\mathbf{x}}_{c,\theta}(\mathbf{z}_{t'}^w))$  $\tilde{\mathbf{x}}_c^w = \frac{\mathbf{z}_{c,t''}^w - (\sigma_{t''}/\sigma_t) \mathbf{z}_t}{\alpha_{t''} - (\sigma_{t''}/\sigma_t) \alpha_t}$   $\triangleright$  Conditional teacher  $\hat{\mathbf{x}}$ 

target

 $\mathbf{z}_{t''}^w = \alpha_{t''} \hat{\mathbf{x}}_\theta(\mathbf{z}_{t'}^w) + \frac{\sigma_{t''}}{\sigma_{t'}} (\mathbf{z}_{t'}^w - \alpha_{t'} \hat{\mathbf{x}}_\theta(\mathbf{z}_{t'}^w))$  $\tilde{\mathbf{x}}^w = \frac{\mathbf{z}_{t''}^w - (\sigma_{t''}/\sigma_t) \mathbf{z}_t}{\alpha_{t''} - (\sigma_{t''}/\sigma_t) \alpha_t}$   $\triangleright$  Unconditional teacher  $\hat{\mathbf{x}}$ 

target

 $\lambda_t = \log[\alpha_t^2/\sigma_t^2]$  $L_\eta = \omega(\lambda_t) (\|\tilde{\mathbf{x}}_c^w - \hat{\mathbf{x}}_{c,\eta}(\mathbf{z}_t, w)\|_2^2 + \|\tilde{\mathbf{x}}^w - \hat{\mathbf{x}}_\eta(\mathbf{z}_t, w)\|_2^2)$  $\eta \leftarrow \eta - \gamma \nabla_\eta L_\eta$ **end while** $\theta \leftarrow \eta$   $\triangleright$  Student becomes next teacher $N \leftarrow N/2$   $\triangleright$  Halve number of sampling steps**end for**

---

Guidance $w$	Number of step	FID ( $\downarrow$ )	IS ( $\uparrow$ )
$w = 0.0$	$1 \times 2$	212.20	3.66
	$16 \times 2$	42.02	7.95
	$64 \times 2$	35.37	8.47
	$128 \times 2$	29.74	8.87
	$256 \times 2$	20.14	9.50
$w = 0.3$	$1 \times 2$	213.07	3.62
	$16 \times 2$	48.74	7.70
	$128 \times 2$	34.28	8.57
	$256 \times 2$	24.54	9.21
$w = 1.0$	$1 \times 2$	214.88	3.54
	$16 \times 2$	64.92	7.21
	$64 \times 2$	48.54	7.62
	$128 \times 2$	42.56	8.00
	$256 \times 2$	32.20	8.81
$w = 2.0$	$1 \times 2$	217.37	3.48
	$16 \times 2$	87.19	6.50
	$64 \times 2$	57.15	7.22
	$128 \times 2$	50.30	7.53
	$256 \times 2$	39.76	8.26
$w = 4.0$	$1 \times 2$	220.11	3.45
	$16 \times 2$	115.57	6.16
	$64 \times 2$	71.45	6.78
	$128 \times 2$	61.75	7.02
	$256 \times 2$	49.21	7.69

Table 5. Distillation results on CIFAR-10 using the naive approach mentioned in Appendix B.8. Note that the naive approach still requires evaluating both a conditional and an unconditional model at each denoising step, and thus requires  $\times 2$  more steps or peak memory than our method. From the evaluated FID/IS scores, we observe that the naive distillation approach is not able to achieve strong performance.

## B.8. Naive distillation approach

A natural approach to progressively distill [33] a classifier-free guided model is to use a distilled student model that follows the same structure as the teacher—that is with a jointly trained distilled conditional and unconditional diffusion component. Denote the pre-trained teacher model  $[\hat{x}_{c,\theta}, \hat{x}_{\theta}]$  and the student model  $[\hat{x}_{c,\eta}, \hat{x}_{\eta}]$ , we provide the training algorithm in Algorithm 5. To sample from the trained model, we can use DDIM deterministic sampler [38] or the proposed stochastic sampler. We follow the training setting in Appendix B.3, use a  $w$ -conditioned model and train the model to condition on the guidance strength  $[0, 4]$ . We observe that the model distilled with Algorithm 5 is not able to generate reasonable samples when the number of sampling is small. We provide the generated samples on CIFAR-10 with DDIM sampler in Fig. 23, and the FID/IS scores in Tab. 5.



Figure 23. Samples using the distillation algorithm mentioned in Appendix B.8. The model is trained with guidance strength  $w \in [0, 4]$  on CIFAR-10. The samples are generated with DDIM (deterministic) sampler at  $w = 0$ . We observe clear artifacts when the number of sampling step is small.

## C. Latent-space distillation

### C.1. Class-conditional generation

#### C.1.1 Training details

In this experiment, we consider class-conditional generation on ImageNet  $256 \times 256$ . We first fine-tune the original  $\epsilon$ -prediction model to a  $v$ -prediction model, and then start from the DDIM teacher model with 512 sampling steps, where we use the output as the target to train our distilled model. For

stage-one, we train the model for 2000 gradient updates with constant loss [10,33]. For stage-two, we train the model with 2000 gradient updates except when the sampling size equals to 1,2, or 4, where we train for 20000 gradient updates. We train the second stage model with SNR-truncation loss [10,33]. For both stages, we train with extra 500 learning rate warm-up steps, where we linearly increase the learning rate from zero to the target learning rate. We use a batch size of 2048 and uniformly sample the guidance strength  $w \in [w_{min} = 0, w_{max} = 14]$  during training.

**Additional results** We provide quantitative results evaluated by precision and recall in Fig. 25. These results confirm a significant performance boost of our method in the small-step regime, especially for 1-4 sampling steps. Our distilled latent diffusion model for 2- and 4-step sampling nearly matches DDIM performance at 32 steps in terms of precision and significantly outperforms it in terms of recall for low numbers of steps. For more qualitative results, see Fig. 25, where we depict random samples for the 1- and 2-step model and contrast them to DDIM sampling.

### C.2. Text-guided image generation

#### C.2.1 Training details

We consider the LAION-5B datasets with resolution  $256 \times 256$  and  $512 \times 512$  in this experiment.

**LAION-5B  $256 \times 256$**  Similar to Appendix C.1, we first fine-tune the original  $\epsilon$ -prediction model to a  $v$ -prediction model. We start from the DDIM teacher model with 512 sampling steps, and use the output as the target to train our distilled model. For stage-one, we train the model for 2000-5000 gradient updates with constant loss [10,33]. For stage-two, we train the model with 2000-5000 gradient updates except when the sampling size equals to 1,2, or 4, where we train for 10000-50000 gradient updates. We train the second stage model with SNR-truncation loss [10,33]. For both stages, we train with extra 100-1000 learning rate warm-up steps, where we linearly increase the learning rate from zero to the target learning rate. We use a batch size of 1024 and uniformly sample the guidance strength  $w \in [w_{min} = 2, w_{max} = 14]$  during training.

Fig. 26 provides a convergence analysis of the different training setting described above. We observe that our method approaches DDIM sampling of the base model after a few thousand training iterations and outperforms it quickly in the 1- and 2-step regime. However, for maximum performance, longer training is required.

**LAION-5B  $512 \times 512$**  Similarly, we first fine-tune the original  $\epsilon$ -prediction model to a  $v$ -prediction model. We start from the DDIM teacher model with 512 sampling steps, and



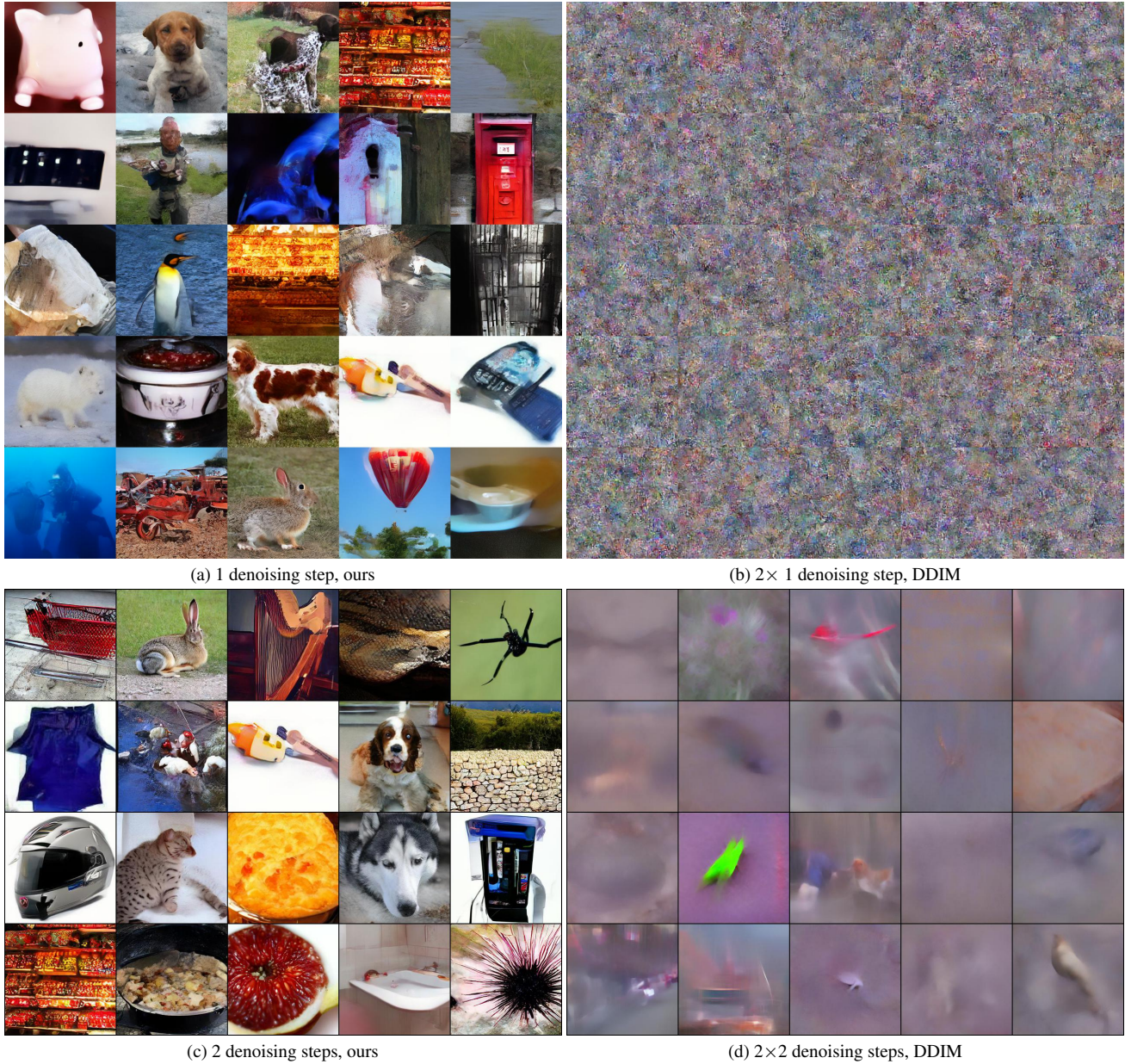


Figure 24. Random  $256 \times 256$  class-conditional samples from our distilled model and from the DDIM teacher for 1 and 2 denoising steps for  $w = 3.0$ .

use the output as the target to train our distilled model. For stage-one, we train the model for 2000-5000 gradient updates with constant loss [10, 33]. For stage-two, we train the model with 2000-5000 gradient updates except when the sampling step equals to 1, 2, or 4, where we train for 10000-50000 gradient updates. We train the second-stage model with SNR-truncation loss [10, 33]. For both stages, we train with extra 1000 learning rate warm-up steps, where we linearly increase the learning rate from zero to the target learning rate. We use a batch size of 512 and uniformly

sample the guidance strength  $w \in [w_{min} = 2, w_{max} = 14]$  during training.

**Additional results** Besides DDIM, we also compare our method here with DPM++-Solver [16, 18], a state-of-the-art sampler that requires no additional training and has achieved good results for  $\geq 10$  sampling steps for latent diffusion models. Unlike our distilled model, this method, similar to DDIM, must use classifier-free guidance to achieve good results. This doubles the number of U-Net evaluations com-

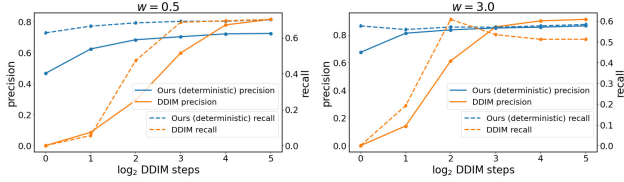


Figure 25. Precision and recall [13] for class-conditional image generation on ImageNet ( $256 \times 256$ ) with distilled latent diffusion. The results are evaluated on 5000 samples. Our distilled latent diffusion model for 2- and 4-step sampling nearly matches DDIM performance at  $32 \times 2$  steps in terms of precision, and strictly outperforms it in terms of recall.

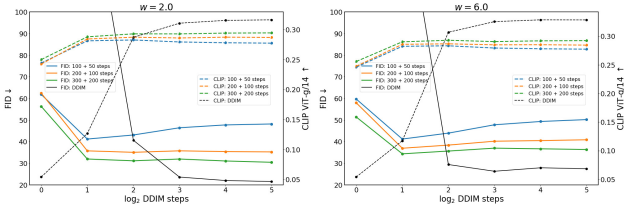


Figure 26. FID and Inception Score for text-guided image generation on LAION ( $256 \times 256$ ) with distilled latent diffusion. The results are evaluated on 5000 captions from COCO2017. We observe that our distillation method approaches DDIM sampling after only a few thousand training steps, see Appendix C.2.

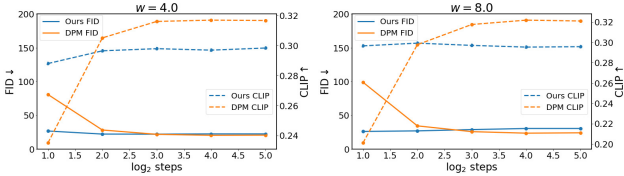


Figure 27. FID and CLIP ViT-g/14 score for text-to-image generation at  $512 \times 512$  px using the distilled *Stable Diffusion* model. The results are evaluated on 5000 captions from the COCO2017 [14] validation set. Our distilled model outperforms the state-of-the-art accelerated sampler *DPM-Solver* (DPM++) [16, 18] in the 2- and 4- step regime. We believe the difference in CLIP scores for  $> 10$ -step sampling can be closed by longer training. We stress that *DPM-Solver*, as DDIM, uses classifier-free guidance during sampling, which requires evaluating both an unconditional and a conditional diffusion model at each denoising step, giving rise to an extra  $\times 2$  overhead compared to our method.

Setting	vs. DDIM (FID)	vs. DPM++ (FID)
2-step, $w = 2.0$	+89.8%	+69.4%
4-step, $w = 2.0$	+68.9%	+32.5%
2-step, $w = 8.0$	+89.5%	+73.7%
4-step, $w = 8.0$	+42.6%	+21.6%

Table 6. Relative performance of our distilled  $512 \times 512$  LAION model compared to DDIM [38] and DPM++ [18] sampling of the base model. Note that DDIM and DPM-Solver use  $2 \times$  more steps than the one listed under “Setting”, as they rely on classifier-free guidance instead of  $w$ -conditioning. This requires DDIM and DPM-Solver to evaluate both an unconditional and a conditional diffusion model at each denoising step, giving rise to the  $\times 2$  overhead.

Setting	vs. DDIM (CLIP)	vs. DPM (CLIP)
2-step, $w = 2.0$	+550%	+27.9%
4-step, $w = 2.0$	+19.2%	+0.1%
2-step, $w = 8.0$	+348%	+47.5%
4-step, $w = 8.0$	+8.6%	+0.6%

Table 7. Relative performance of our distilled  $512 \times 512$  LAION model compared to DDIM [38] and DPM-Solver (DPM++) [16, 18] sampling of the base model. Note that DDIM and DPM use  $2 \times$  more steps than the one listed under “Setting”, as they rely on classifier-free guidance instead of  $w$ -conditioning. This requires DDIM and DPM to evaluate both an unconditional and a conditional diffusion model at each denoising step, giving rise to the  $\times 2$  overhead. We use CLIP ViT-g/14 for evaluation [7, 25].

pared to our  $w$ -conditional approach.

We provide a qualitative comparison of these sampling methods in Fig. 28, where we clearly see the benefits of our distillation approach for low numbers of sampling steps: our method produces sharper and more coherent results than the training-free samplers. This behavior is reflected by the quantitative FID and CLIP analysis in Fig. 27 and Tab. 6, Tab. 7. While the speed-up here is not quite as significant as in pixel-space, our method still achieves very good results with 2 or 4 sampling steps. Our approach further reduces the maximum memory or denoising step by a half compared to existing methods due to  $w$ -conditioning (since here we no longer need to evaluate both the unconditional model and conditional model for classifier-free guidance, we only need one distilled  $w$ -conditional model). We hope that our work will lead to progress in real-time applications of general high-resolution text-to-image systems.

We also provide human evaluation results by leveraging Amazon Mechanical Turk. We generate images using text prompts from [45]. We compare our distilled model sampled using 2 or 4 denoising steps with DDIM and DPM++ solver sampled using  $2 \times 2$  or  $4 \times 2$  denoising steps. For each setting, we generate 100 HITs each with 17 pair-wise comparisons between samples generated with our approach and the baseline. In each of the question, the user is shown the text prompt used to generate the image and asked to select the image that looks better to them. We provide a snapshot of our user interface in Fig. 29. We provide the results in Tab. 9. Although we observe noisy answers (for instance some user would prefer the right image to the left image in Fig. 29c), our distilled model still consistently outperforms the baselines in all the settings we considered in Tab. 9. To get higher-quality user feedback and reduce the noise in the answers, in the future work, we will perform a new human evaluation with a larger sample size and extra constraints to ensure the quality of the response. We will also build a framework to automatically ignore HITs with random selections.



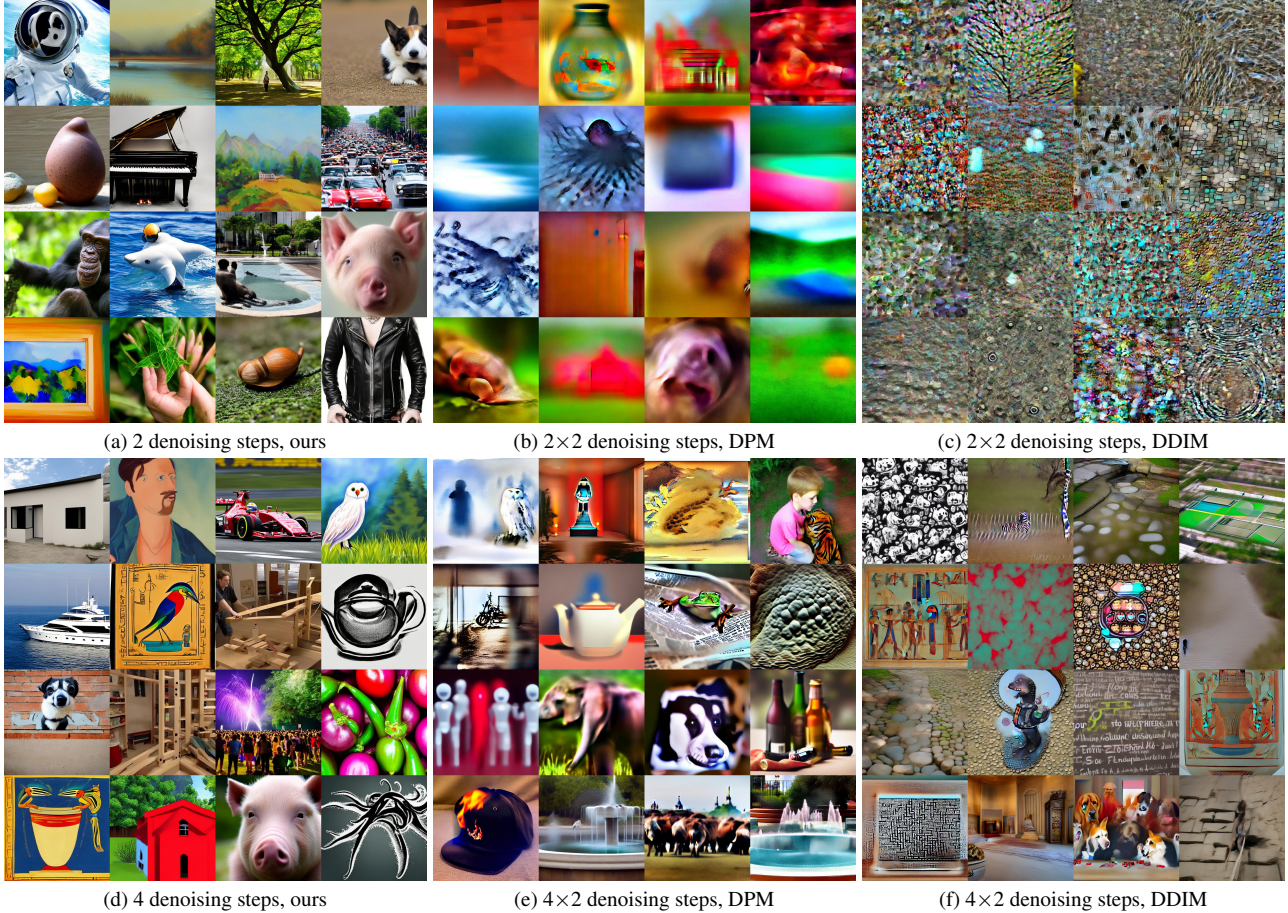


Figure 28. Random  $512 \times 512$  text-guided samples from our distilled *Stable Diffusion* model compared to the DDIM teacher and DPM-solver for 2 and 4 denoising steps for  $w = 11.5$ .

### C.3. Text-guided image-to-image translation

#### C.3.1 Training details

We use the model trained for text-guided image generation. The training details can be found in Appendix C.2.

#### C.3.2 Extra analysis

We provide more analysis on the trade-off between sample quality, controllability and efficiency in Fig. 30 and Fig. 31. Similar to [20], we also observe a trade-off between realism, controllability and faithfulness as we increase the initial perturbed noise level: the more noise we add, the more aligned the images are to the text prompt, but less faithful to the input image (see Fig. 30 and Fig. 31).

### C.4. Image inpainting

#### C.4.1 Training details

Similar to our previous experiments, we fine-tune the  $\epsilon$ -prediction model to a  $v$ -prediction model, using the large

Setting	Ours (FID ↓)	DDIM (FID ↓)
2-step, $w = 4.0$	29.50	109.35
4-step, $w = 4.0$	24.90	26.89
2-step, $w = 11.0$	31.43	105.71
4-step, $w = 11.0$	24.36	27.22

Table 8. Quantitative inpainting results as evaluated by FID. We evaluate on 2000 examples from COCO2017. Note that DDIM, which is evaluated with classifier-free guidance, uses two times more function evaluations than the one listed under “Setting”.

mask generation scheme suggested in LAMA [42] and train on LAION-5B at  $512 \times 512$  resolution. We start from the DDIM teacher model with 512 sampling steps, and use the output as the target to train our distilled model. For stage-one, we train the model for 2000 gradient updates with constant loss [10, 33]. For stage-two, we train the model with 10000 gradient updates except when the sampling size equals to 1 or 2, where we train for 5000 gradient updates. We train the second stage model with SNR-truncation loss [10, 33]. For



**About this HIT:**

- Please only participate in this HIT if you have normal color vision.
- It should take about 1 minute.
- You will take part in an experiment involving visual perception. You'll see a text prompt and a series of pairs of images. In each pair, given the text prompt, the images are "fake" images generated using a computer program. Choose the image that looks **more reasonable** to you. Your selection should be based on how **realistic** and **less blurry** the image is, and whether the image **follows the text prompt**.

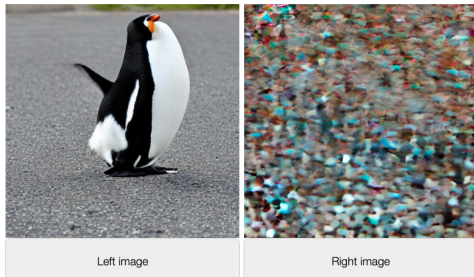
Start!

Given the text prompt: "a glass of orange juice", how would you imagine this image to look like?  
Choose the image that looks **more reasonable** to you.  
Your selection should be based on how **realistic** and **less blurry** the image is, and whether it **follows the text prompt**.

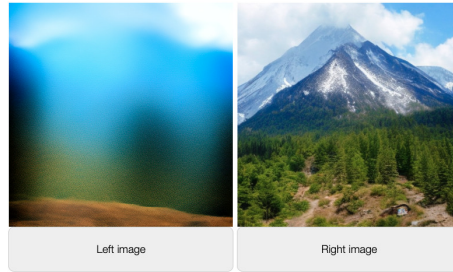


(a) Instructions for the human evaluators on Amazon Mechanical Turk. (b) Images generated by our 4-step distillation model (left) and images generated by the  $4 \times 2$ -step baseline (right).

Given the text prompt: "a penguin standing on a sidewalk", how would you imagine this image to look like?  
Choose the image that looks **more reasonable** to you.  
Your selection should be based on how **realistic** and **less blurry** the image is, and whether it **follows the text prompt**.



Given the text prompt: "a mountain", how would you imagine this image to look like?  
Choose the image that looks **more reasonable** to you.  
Your selection should be based on how **realistic** and **less blurry** the image is, and whether it **follows the text prompt**.



(c) Images generated by our 2-step distillation model (left) and (d) Images generated by the  $2 \times 2$ -step baseline (left) and images generated by the  $2 \times 2$ -step baseline (right).

Figure 29. A snapshot of the human evaluation interface we used on Amazon Mechanical Turk.

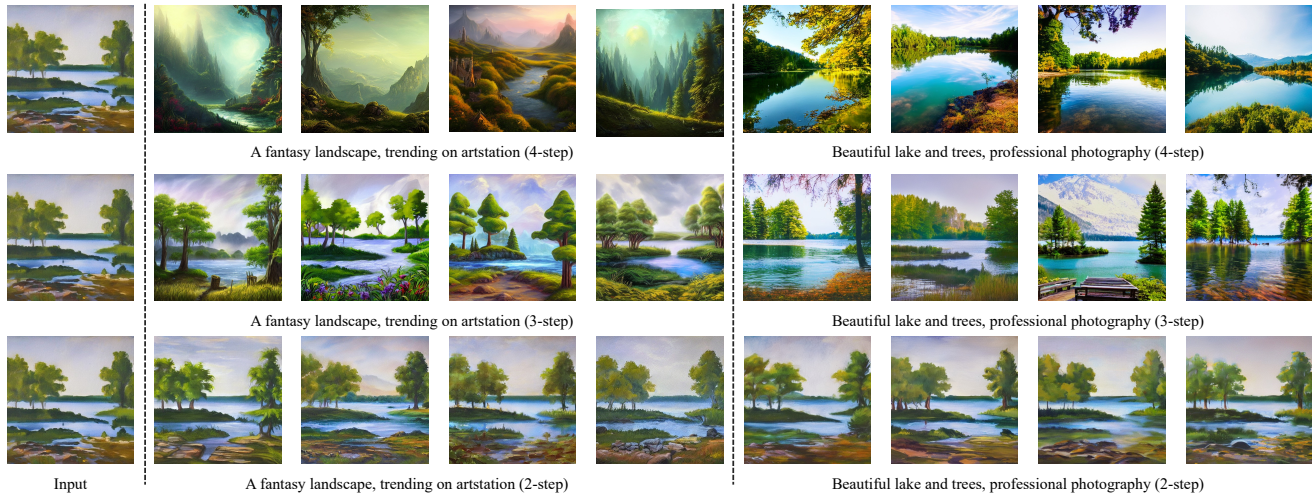


Figure 30. In this example, we study the trade-off between efficiency, realism, and controllability for guided image translation with SDEdit [20]. We use a 4-step distilled text-guided image generation model trained on LAION-5B ( $512 \times 512$ ). The training detail is discussed in Appendix C.2. Given an input image (guide), we consider perturbing the input image with different noise level, with 2 denoising step corresponding to perturb the image with around 50% noise, and 4 denoising step corresponding to perturb the image with around 100% noise according to the DDIM noise schedule. We observe that the more noise we perturb, the more aligned the images are with the text prompt, but the less faithful they are to the input image.



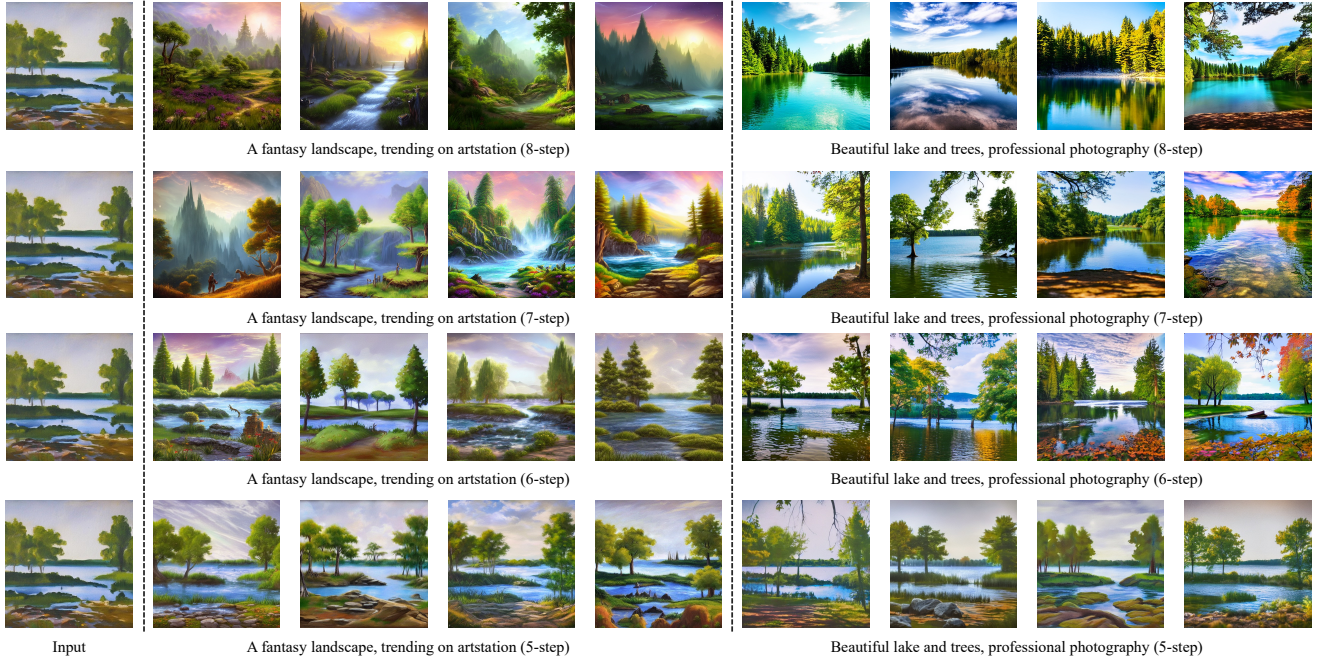


Figure 31. In this example, we study the trade-off between efficiency, realism, and controllability for guided image translation with SDEdit [20]. We use a 8-step distilled text-guided image generation model trained on LAION-5B ( $512 \times 512$ ). The training detail is discussed in Appendix C.2. Given an input image (guide), we consider perturbing the input image with different noise level, with 5 denoising step corresponding to perturb the image with around 60% noise, and 8 denoising step corresponding to perturb the image with around 100% noise according to the DDIM noise schedule. We observe that the more noise we perturb, the more aligned the images are with the text prompt, but the less faithful they are to the input image.



Figure 32. Random  $512 \times 512$  inpainting samples from our distilled model and from the DDIM teacher for 2 denoising steps for  $w = 11.0$ .

both stages, we train with extra 1000 learning rate warm-up steps, where we linearly increase the learning rate from zero to the target learning rate. We use a batch size of 512 and uniformly sample the guidance strength  $w \in [w_{min} = 2, w_{max} = 14]$  during training.

**Additional evaluation results** A quantitative comparison with DDIM sampling at low sampling numbers of sampling steps can be found in Tab. 8, additional samples are in Fig. 32.

Ours	Baseline	Our method is better ( $\uparrow$ )
Distillation 2-step	DDIM 2 $\times$ 2-step	66.32%
Distillation 2-step	DPM++ 2 $\times$ 2-step	68.97%
Distillation 2-step	DDIM 4 $\times$ 2-step	57.44%
Distillation 2-step	DPM++ 4 $\times$ 2-step	59.88%
Distillation 4-step	DDIM 4 $\times$ 2-step	67.36%
Distillation 4-step	DPM++ 4 $\times$ 2-step	64.71%

Table 9. Human evaluation on text-guided image generation. Here the model is trained on LAION-5B (512 $\times$ 512). We leverage Amazon Mechanical Turk for human evaluation. We perform pairwise comparison between our method and the baselines. We compare our method using 2 or 4 denoising steps with DDIM [38] and DPM++ [18] samplers using 2 $\times$ 2 or 4 $\times$ 2 denoising steps. We use a guidance strength of 12.5 for all methods. For each setting, we distribute 100 HITs each with 17 pairwise comparison questions. We show MTurk workers the text prompt as well as the two generated images, and then ask them to select the one they think is better. We provide a snapshot of the interface in Fig. 29. In the table, we report the percentage that the MTurk workers think our method is better than the baseline. Although, we observe noise in the response (some user would prefer the right image to the left image in Fig. 29c), our method still consistently outperform the baselines in all settings. For the future work, we will incorporate schemes to ignore invalid HITs with random answers. We will also perform another human evaluation study with larger sample sizes and more constraints to ensure high-quality responses.

## D. Extra samples for pixel-space distillation

In this section, we provide extra samples for the pixel-space distillation models. We generate samples using the deterministic sampler (see Algorithm 2) and the stochastic sampler (see Algorithm 3).



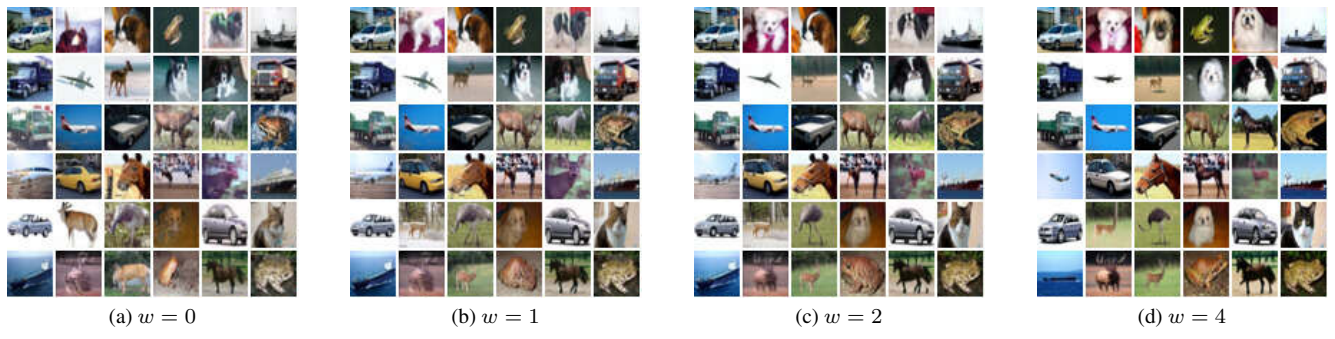


Figure 33. Ours (deterministic in pixel-space) on CIFAR-10. Distilled 256 sampling steps.

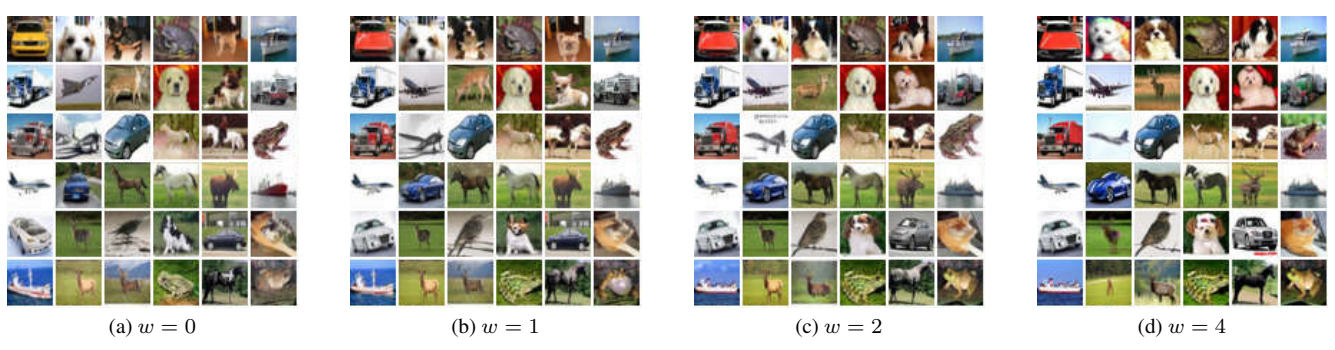


Figure 34. Ours (stochastic in pixel-space) on CIFAR-10. Distilled 256 sampling steps.

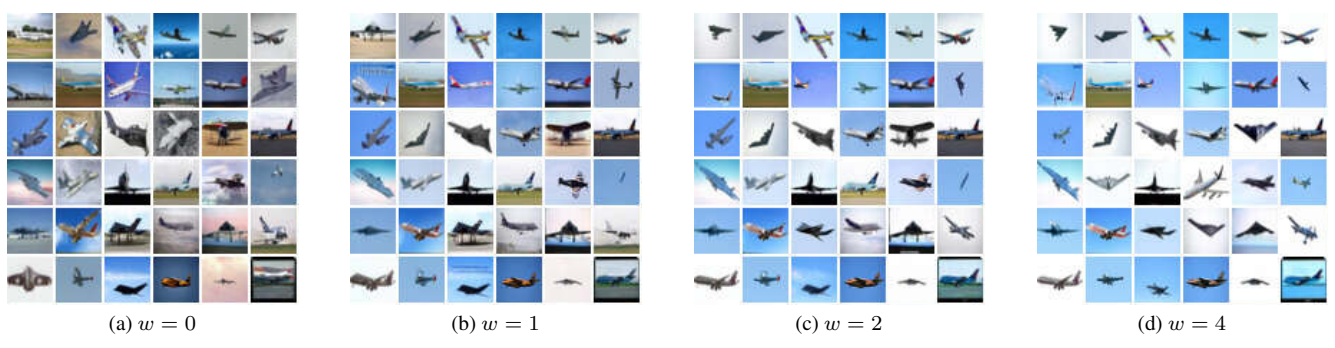


Figure 35. Ours (deterministic in pixel-space) on CIFAR-10. Distilled 256 sampling steps. Class-conditioned samples.

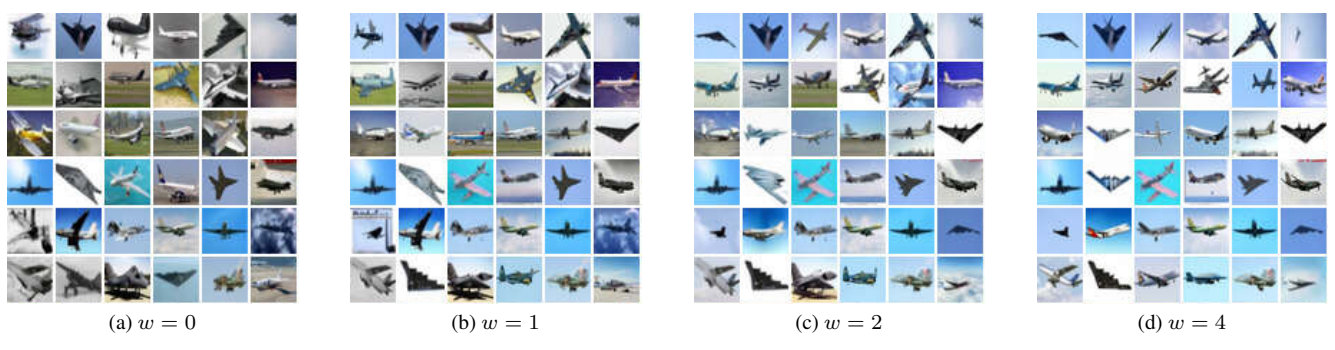


Figure 36. Ours (stochastic in pixel-space) on CIFAR-10. Distilled 256 sampling steps. Class-conditioned samples.



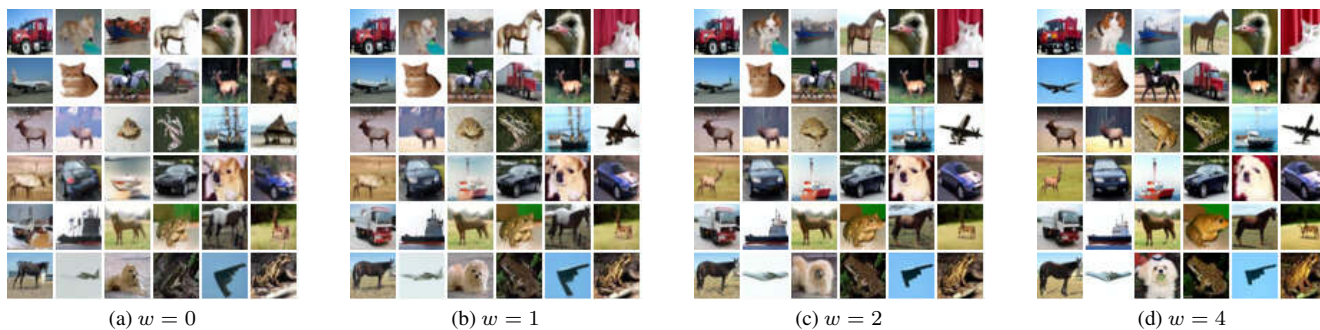


Figure 37. Ours (deterministic in pixel-space) on CIFAR-10. Distilled 4 sampling steps.



Figure 38. Ours (stochastic in pixel-space) on CIFAR-10. Distilled 4 sampling steps.

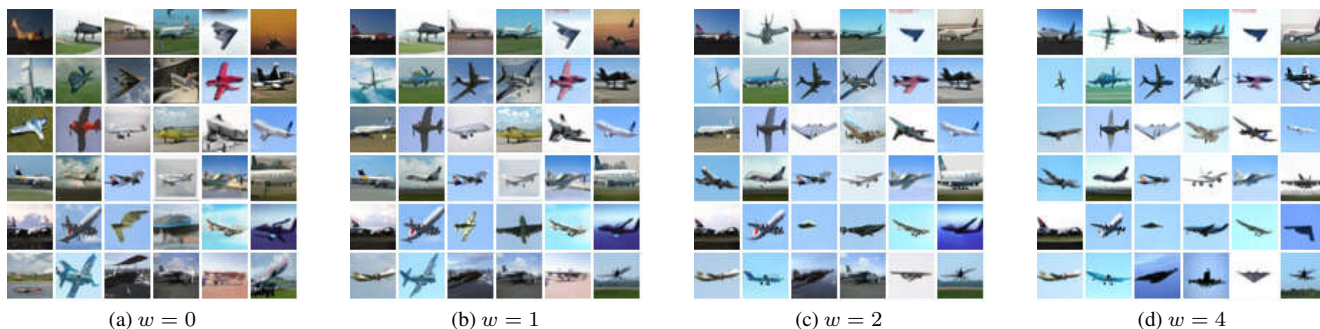


Figure 39. Ours (deterministic in pixel-space) on CIFAR-10. Distilled 4 sampling steps. Class-conditioned samples.

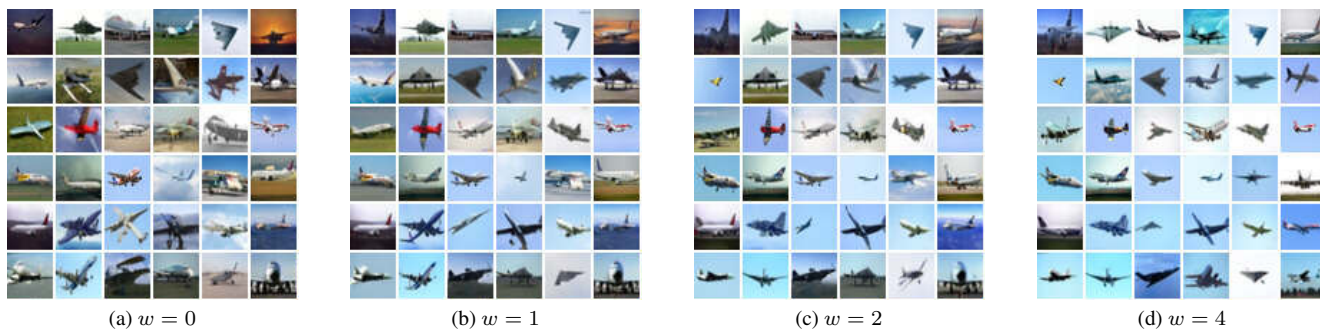


Figure 40. Ours (stochastic in pixel-space) on CIFAR-10. Distilled 4 sampling steps. Class-conditioned samples.



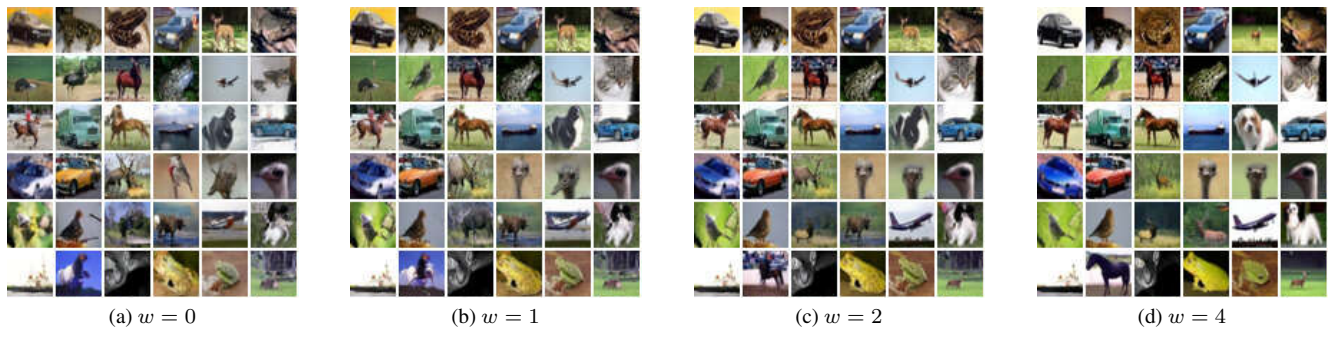


Figure 41. Ours (deterministic in pixel-space) on CIFAR-10. Distilled 2 sampling steps.

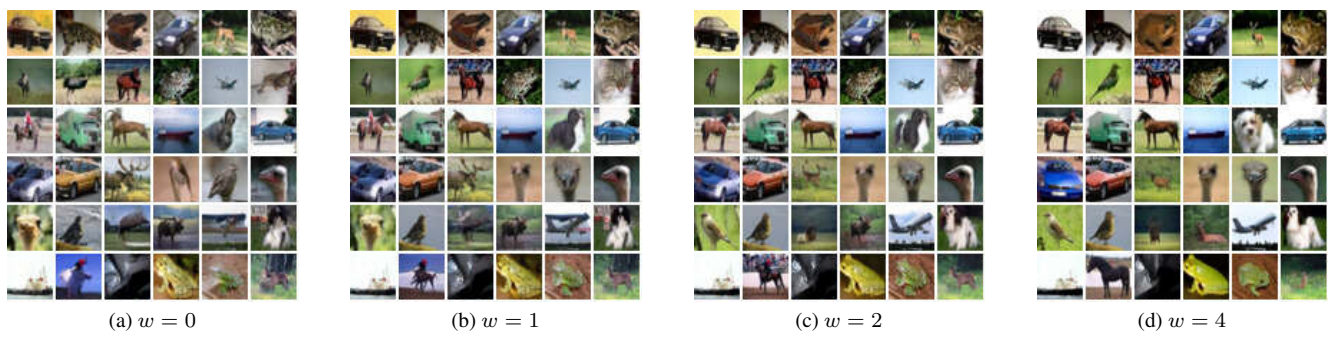


Figure 42. Ours (stochastic in pixel-space) on CIFAR-10. Distilled 2 sampling steps.

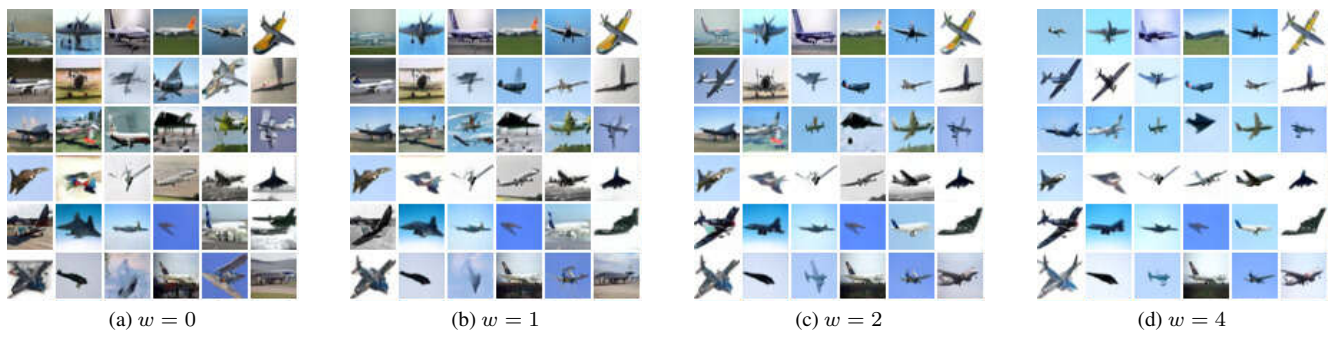


Figure 43. Ours (deterministic in pixel-space) on CIFAR-10. Distilled 2 sampling steps. Class-conditioned samples.

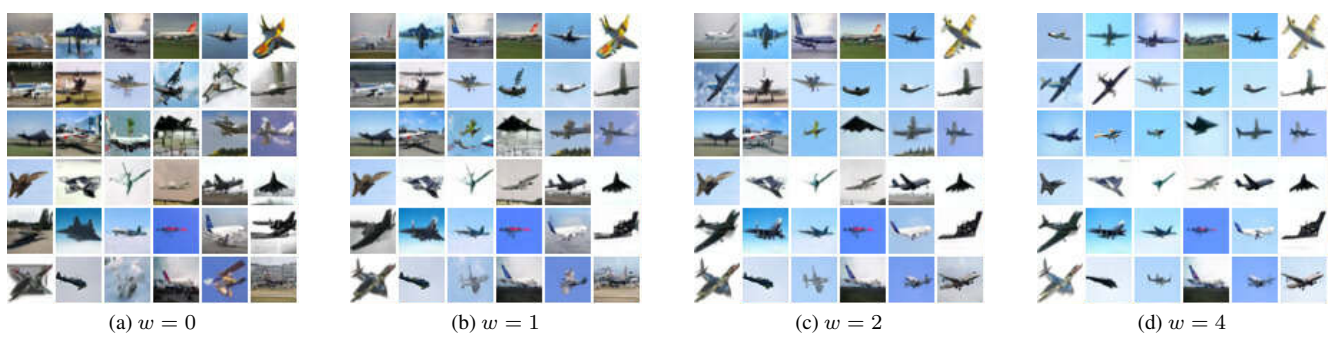


Figure 44. Ours (stochastic in pixel-space) on CIFAR-10. Distilled 2 sampling steps. Class-conditioned samples.



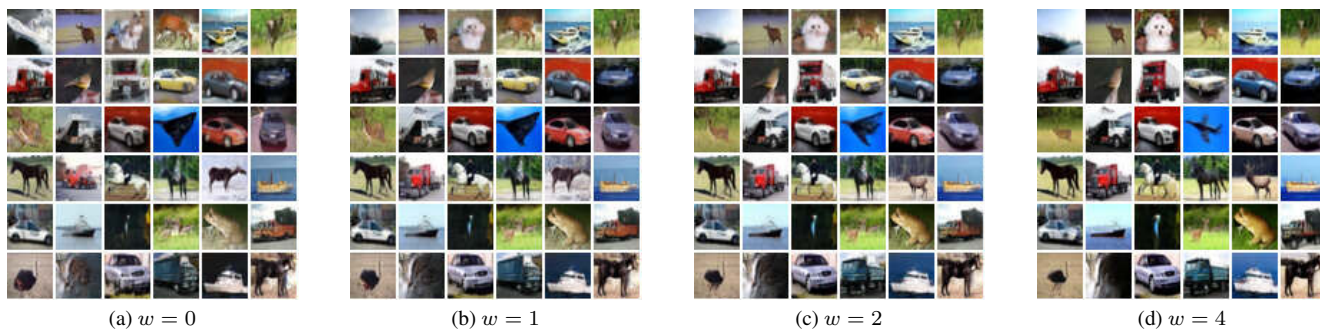


Figure 45. Ours (deterministic in pixel-space) on CIFAR-10. Distilled 1 sampling step.

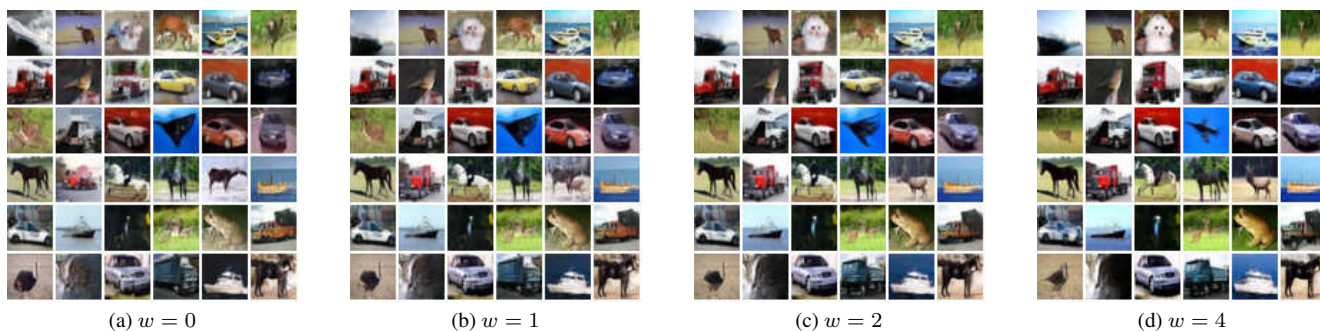


Figure 46. Ours (stochastic in pixel-space) on CIFAR-10. Distilled 1 sampling step.

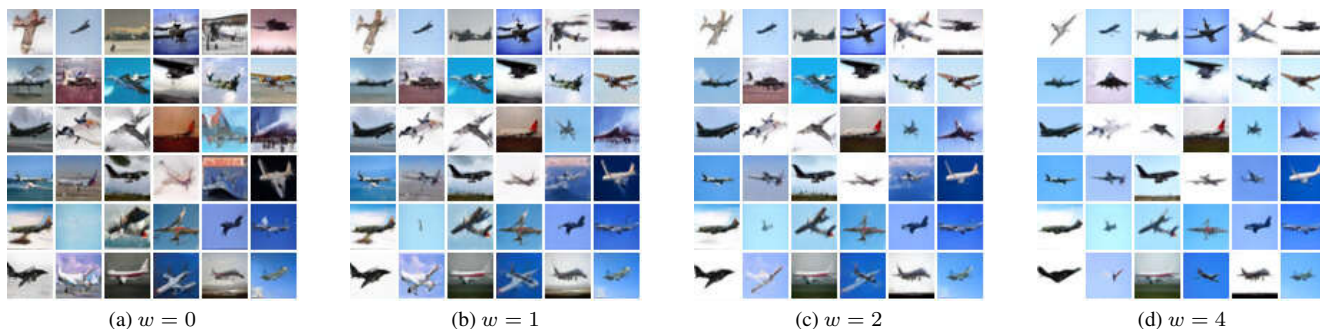


Figure 47. Ours (deterministic in pixel-space) on CIFAR-10. Distilled 1 sampling step. Class-conditioned samples.

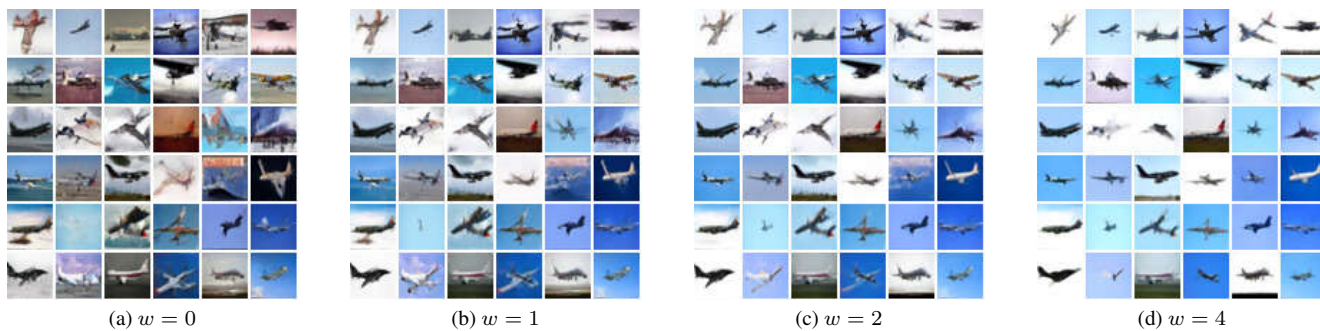


Figure 48. Ours (stochastic in pixel-space) on CIFAR-10. Distilled 1 sampling step. Class-conditioned samples.



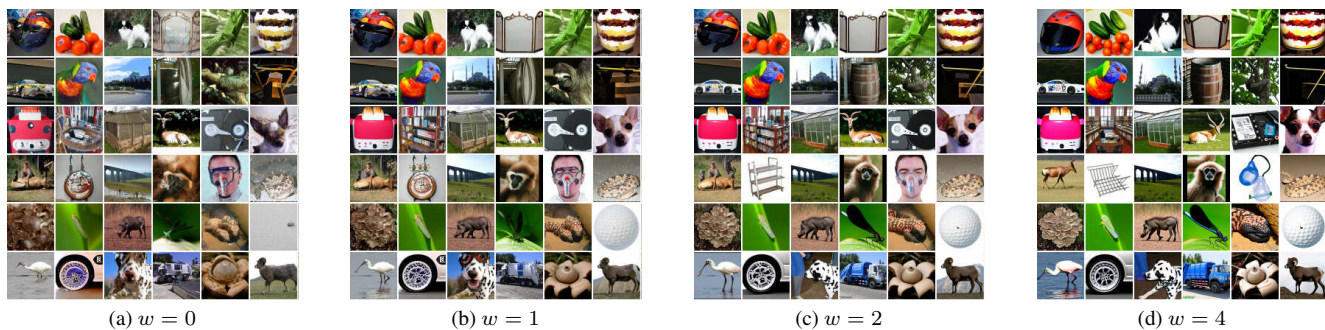


Figure 49. Ours (deterministic in pixel-space) on ImageNet 64x64. Distilled 256 sampling steps.

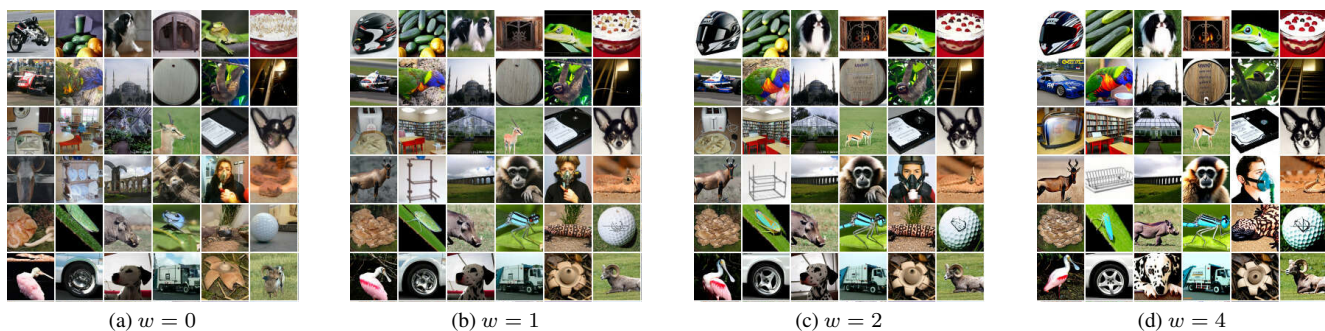


Figure 50. Ours (stochastic in pixel-space) on ImageNet 64x64. Distilled 256 sampling steps.



Figure 51. Ours (deterministic in pixel-space) on ImageNet 64x64. Distilled 256 sampling steps. Class-conditioned samples.

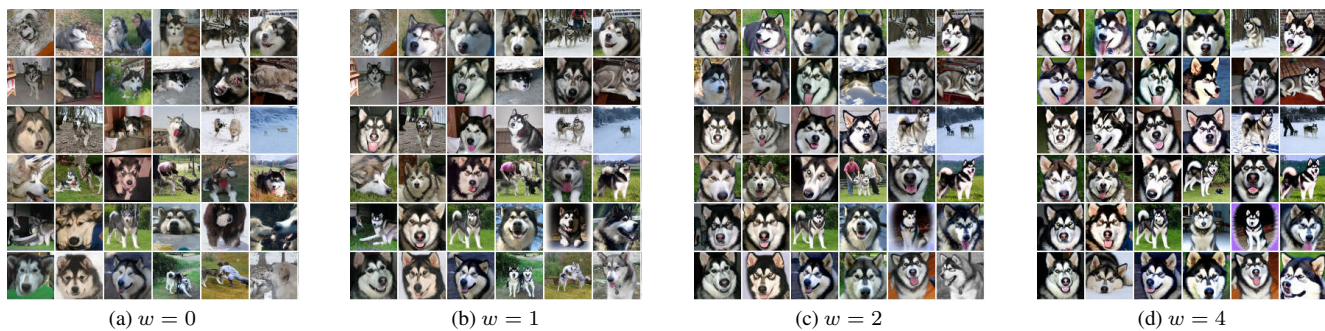


Figure 52. Ours (stochastic in pixel-space) on ImageNet 64x64. Distilled 256 sampling steps. Class-conditioned samples.



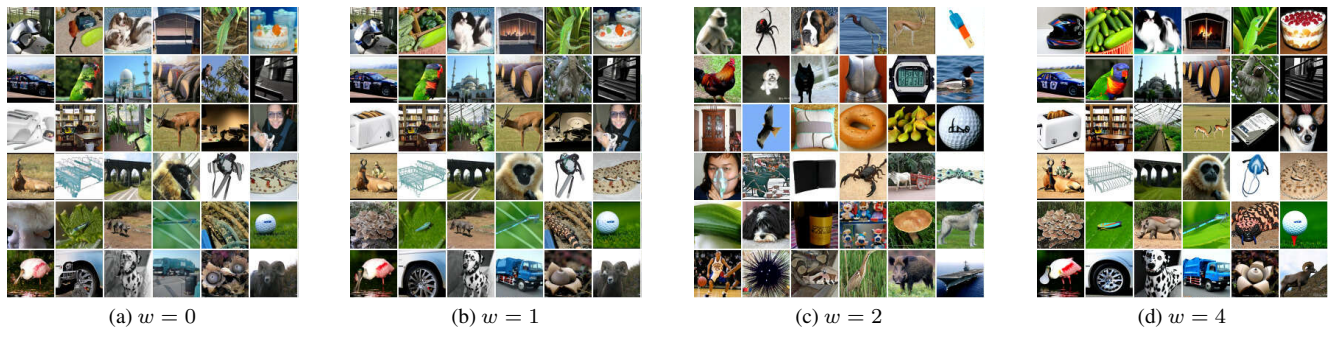


Figure 53. Ours (deterministic in pixel-space) on ImageNet 64x64. Distilled 8 sampling step.

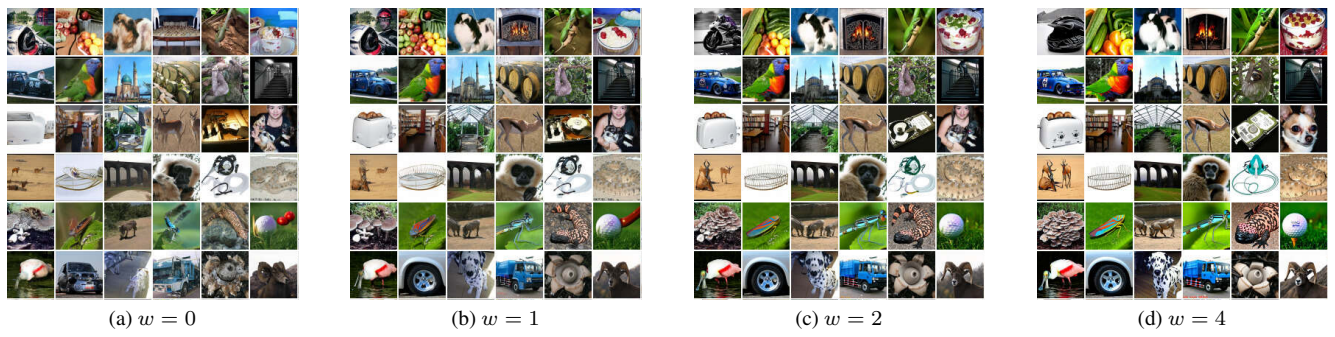


Figure 54. Ours (stochastic in pixel-space) on ImageNet 64x64. Distilled 8 sampling step.

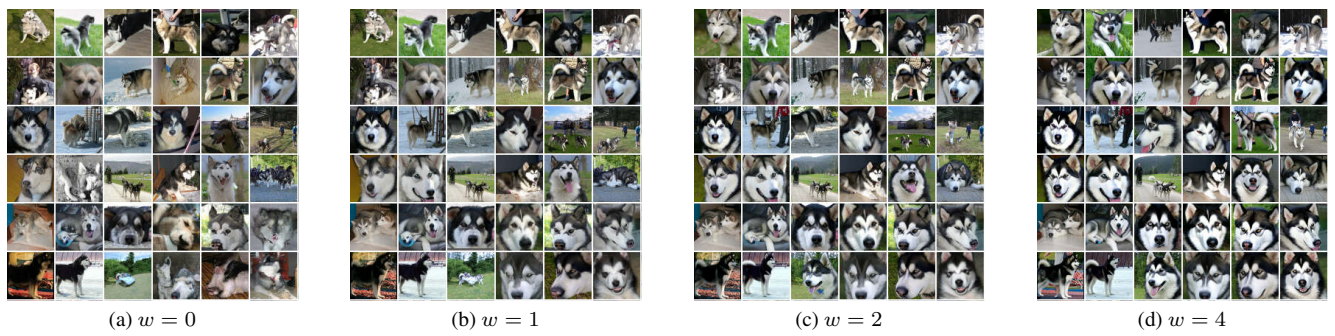


Figure 55. Ours (deterministic in pixel-space) on ImageNet 64x64. Distilled 8 sampling step. Class-conditioned samples.

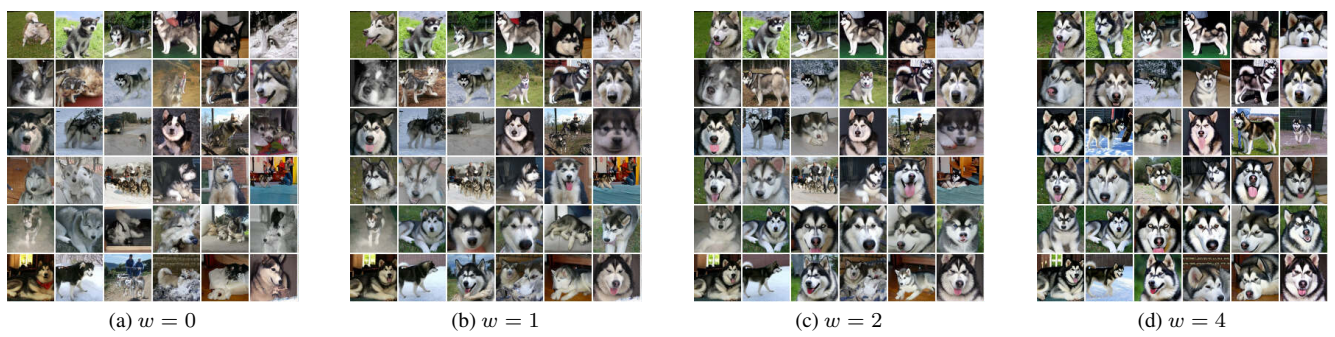


Figure 56. Ours (stochastic in pixel-space). Distilled 8 sampling step. Class-conditioned samples.



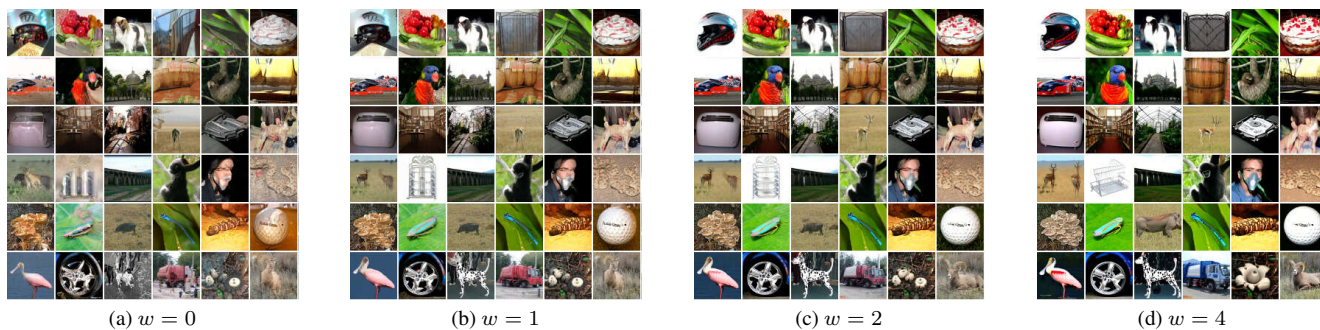


Figure 57. Ours (deterministic in pixel-space) on ImageNet 64x64. Distilled 2 sampling steps.

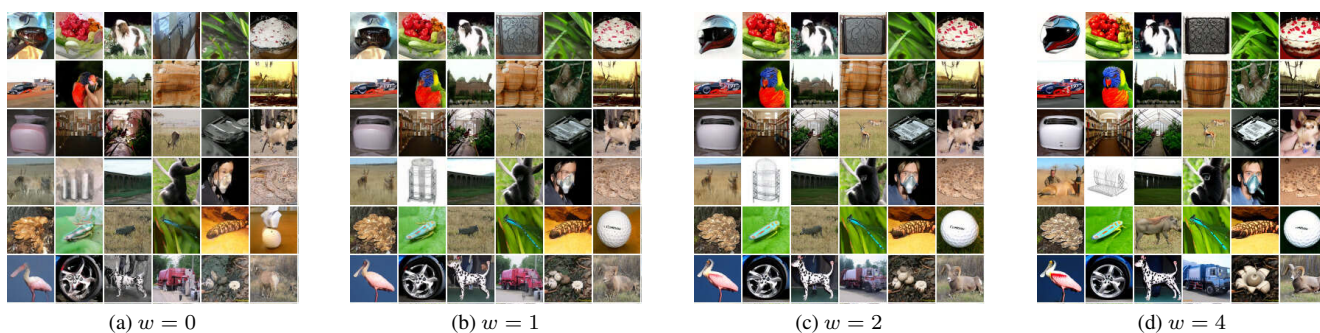


Figure 58. Ours (stochastic in pixel-space) on ImageNet 64x64. Distilled 2 sampling steps.

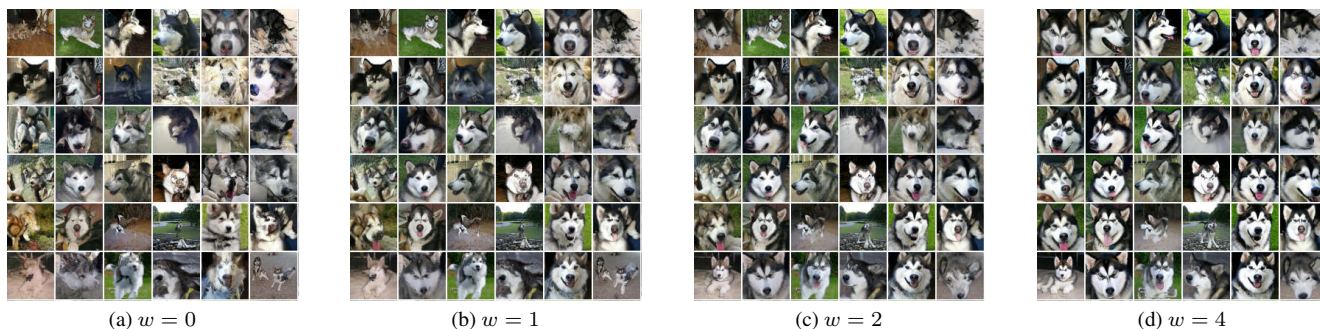


Figure 59. Ours (deterministic in pixel-space) on ImageNet 64x64. Distilled 2 sampling steps. Class-conditioned samples.

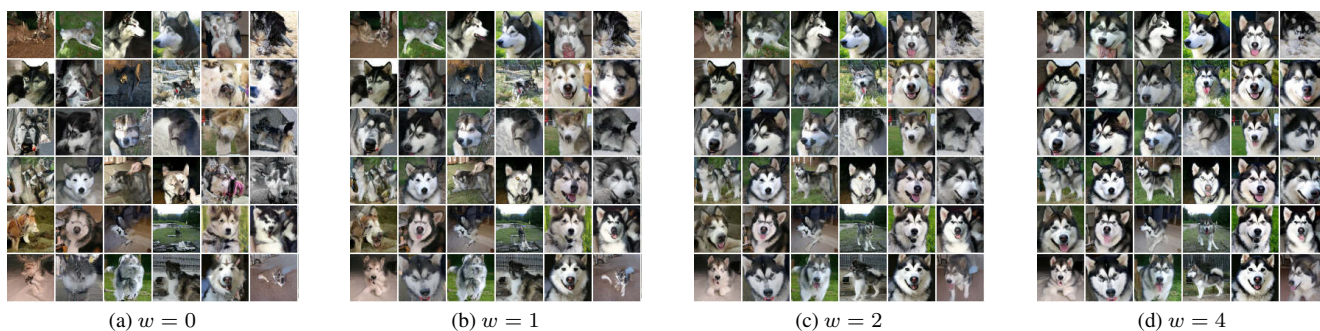


Figure 60. Ours (stochastic in pixel-space) on ImageNet 64x64. Distilled 2 sampling steps. Class-conditioned samples.



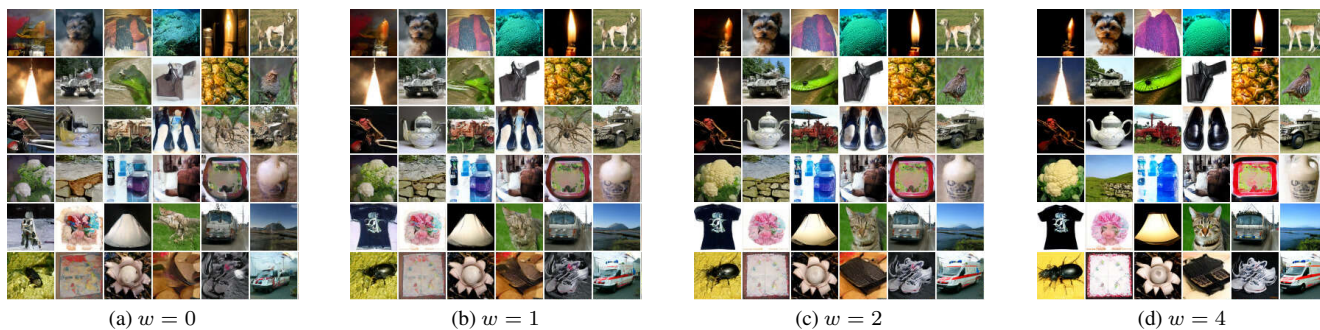


Figure 61. Ours (deterministic in pixel-space) on ImageNet 64x64. Distilled 1 sampling step.

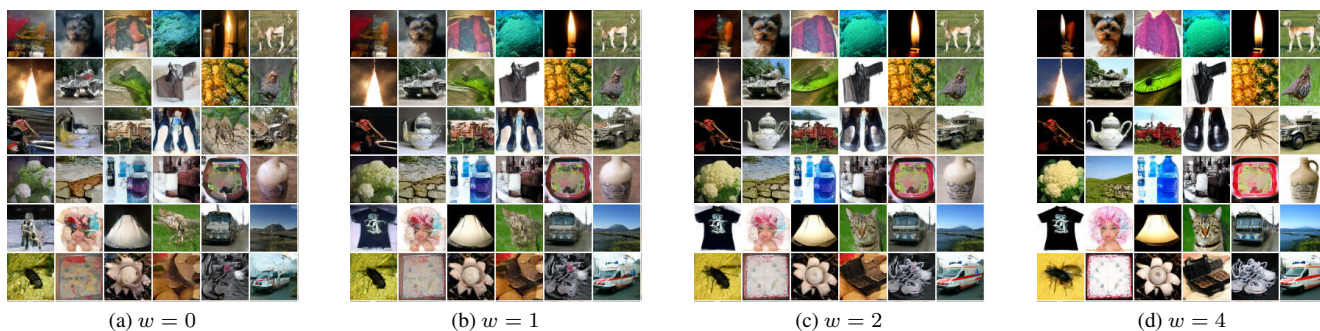


Figure 62. Ours (stochastic in pixel-space) on ImageNet 64x64. Distilled 1 sampling step.

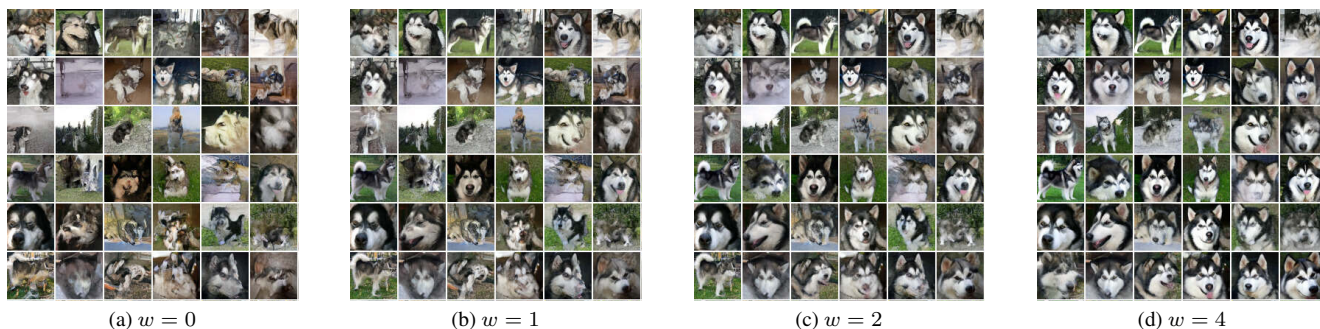


Figure 63. Ours (deterministic in pixel-space) on ImageNet 64x64. Distilled 1 sampling step. Class-conditioned samples.

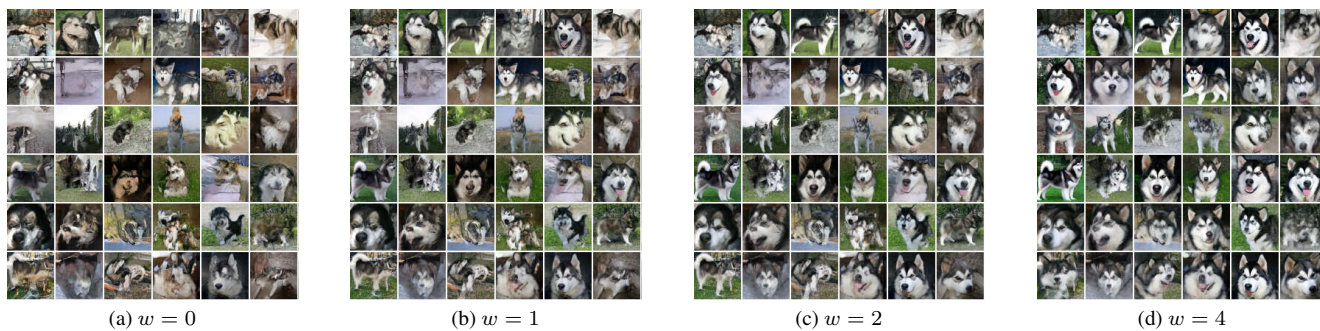


Figure 64. Ours (stochastic in pixel-space) on ImageNet 64x64. Distilled 1 sampling step. Class-conditioned samples.