

# Diffusion Models without Classifier-free Guidance

Zhicong Tang<sup>1</sup> Jianmin Bao<sup>2</sup> Dong Chen<sup>2</sup> Baining Guo<sup>2</sup>

## Abstract

This paper presents Model-guidance (MG), a novel objective for training diffusion model that addresses and removes the commonly used Classifier-free guidance (CFG). Our innovative approach transcends the standard modeling of solely data distribution to incorporating the posterior probability of conditions. The proposed technique originates from the idea of CFG and is easy yet effective, making it a plug-and-play module for existing models. Our method significantly accelerates the training process, doubles the inference speed, and achieve exceptional quality that parallel and even surpass concurrent diffusion models with CFG. Extensive experiments demonstrate the effectiveness, efficiency, scalability on different models and datasets. Finally, we establish state-of-the-art performance on ImageNet 256 benchmarks with an FID of 1.34. Our code is available at [github.com/tzco/Diffusion-wo-CFG](https://github.com/tzco/Diffusion-wo-CFG).

## 1. Introduction

Diffusion models (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020; Song et al., 2021a;b) have become the cornerstone of many successful generative models, *e.g.* image generation (Dhariwal & Nichol, 2021; Nichol et al., 2022; Rombach et al., 2022; Podell et al., 2024; Chen et al., 2024) and video generation (Ho et al., 2022; Blattmann et al., 2023; Gupta et al., 2025; Polyak et al., 2024; Wang et al., 2024) tasks. However, diffusion models also struggle to generate “low temperature” samples (Ho & Salimans, 2021; Karras et al., 2024) due to the nature of training objectives, and techniques such as Classifier guidance (Dhariwal & Nichol, 2021) and Classifier-free guidance (CFG) (Ho & Salimans, 2021) are proposed to improve performances.

Despite its advantage and ubiquity, CFG has several drawbacks (Karras et al., 2024) and poses challenges to effective implementations (Kynkäänniemi et al., 2024) of diffusion

<sup>1</sup>Tsinghua University <sup>2</sup>Microsoft Research Asia.

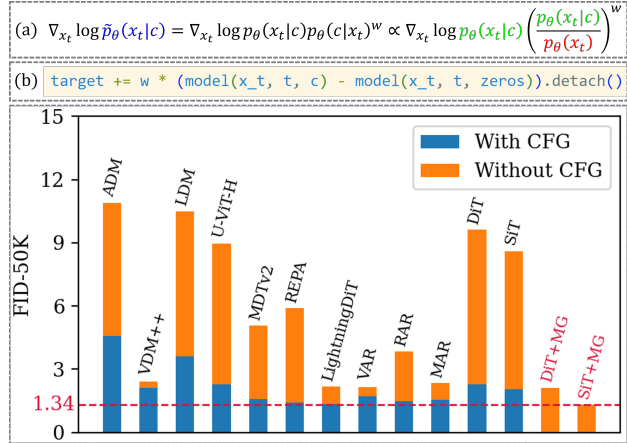


Figure 1: **We propose Model-guidance (MG), removing Classifier-free guidance (CFG) for diffusion models and achieving state-of-the-art on ImageNet with FID of 1.34.** (a) Instead of running models twice during inference (green and red), MG directly learns the final distribution (blue). (b) MG requires only one line of code modification while providing excellent improvements. (c) Comparing to concurrent methods, MG yields lowest FID even without CFG.

models. One critical limitation is the simultaneous training of unconditional model apart from the main diffusion model. The unconditional model is typically implemented by randomly dropping the condition of training pairs and replacing with an manually defined empty label. The introduction of additional tasks may reduce network capabilities and lead to skewed sampling distributions (Karras et al., 2024; Kynkäänniemi et al., 2024). Furthermore, CFG requires two forward passes per denoising step during inference, one for the conditioned and another for the unconditioned model, thereby significantly escalating the computational costs.

In this work, we propose Model-guidance (MG), an innovative method for diffusion models to effectively circumvent CFG and boost performances, thereby eliminating the limitations above. We propose a novel objective that transcends from simply modeling the data distribution to incorporating the posterior probability of conditions. Specifically, we leverage the model itself as an implicit classifier and directly learn the score of calibrated distribution during training.

As depicted in Figure 1, our proposed method confers mul-

multiple substantial breakthroughs. First, it significantly refines generation quality and accelerates training processes, with experiments showcasing a  $\geq 6.5\times$  convergence speedup than vanilla diffusion models with excellent quality. Second, the inference speed is doubled with our method, as each denoising step needs only one network forward in contrast to two in CFG. Besides, it is easy to implement and requires only one line of code modification, making it a plug-and-play module of existing diffusion models with instant improvements. Finally, it is an end-to-end method that excels traditional two-stage distillation-based approaches and even outperforms CFG in generation performances.

We conduct comprehensive experiments on the prevalent Imagenet (Deng et al., 2009; Russakovsky et al., 2015) benchmarks with  $256 \times 256$  and  $512 \times 512$  resolution and compare with a wide variates of concurrent models to attest the effectiveness of our proposed method. The evaluation results demonstrate that our method not only parallels and even outperforms other approaches with CFG, but also scales to different models and datasets, making it a promising enhancement for diffusion models. In conclusion, we make the following contribution in this work:

- We proposed a novel and effective method, Model-guidance (MG), for training diffusion models.
- MG removes CFG for diffusion models and greatly accelerates both training and inference process.
- Extensive experiments with SOTA results on ImageNet demonstrate the usefulness and advantages of MG.

## 2. Background

### 2.1. Diffusion and Flow Models

**Diffusion models** (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020; Song et al., 2021a;b) are a class of generative models that utilize forward and reverse stochastic processes to model complex data distributions.

The forward process adds noise and transforms data samples into Gaussian distributions as

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}), \quad (1)$$

where  $x_t$  represents the noised data at timestep  $t$  and  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$  is the noise schedule.

Conversely, the reverse process learns to denoise and finally recover the original data distribution, which aims to reconstruct score (Sohl-Dickstein et al., 2015; Song et al., 2021b) from the noisy samples  $x_t$  by learning

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)), \quad (2)$$

where  $\mu_\theta$  and  $\Sigma_\theta$  are mean and variance and commonly predicted by neural networks.

In common implementations, the training of diffusion mod-

els leverages a re-parameterized objective that directly predicts the noise at each step (Ho et al., 2020)

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{t, x_0, \epsilon} \|\epsilon_\theta(x_t, t) - \epsilon\|^2, \quad (3)$$

where  $x_t$  is derived from the forward process in Equation (1) with  $x_0$  and  $\epsilon$  drawn from dataset and Gaussian noises.

Conditional diffusion models allow users to generate samples aligned with specified demands and precisely control the contents of samples. In this case, the generation process is manipulated with give conditions  $c$ , such as class labels or text prompts, where network functions are  $\epsilon_\theta(x_t, t, c)$ .

**Flow Models** (Lipman et al., 2023; Liu et al., 2023; Albergo et al., 2023; Tong et al., 2024) are another emerging type of generative models similar to diffusion models. Flow models utilize the concept of Ordinary Differential Equations (ODEs) to bridge the source and target distribution and learn the directions from noise pointing to ground-truth data.

The forward process of flow models is defined as an Optimal Transport (OT) interpolant (McCann, 1997)

$$x_t = (1 - t)x_0 + t\epsilon, \quad (4)$$

and the loss function takes the form (Lipman et al., 2023)

$$\mathcal{L}_{\text{FM}} = \mathbb{E}_{t, x_0, \epsilon} \|u_\theta(x_t) - u_t(x_t|x_0)\|^2, \quad (5)$$

where the ground-truth conditional flow is given by

$$u_t(x_t|x_0) = x_0 - \epsilon. \quad (6)$$

### 2.2. Classifier-Free Guidance

Classifier-free guidance (CFG) (Ho & Salimans, 2021) is a widely adopted technique in conditional diffusion models to enhance generation performance and alignment to conditions. It provides an explicit control of the focus on conditioning variables and avoids to sample within the ‘‘low temperature’’ regions with low quality.

The key design of CFG is to combine the posterior probability and utilize Bayes’ rule during inference time. To facilitate this, it is required to train both conditional and unconditional diffusion models. In particular, CFG trains the models to predict

$$\epsilon_\theta(x_t, t, c) \propto -\nabla_{x_t} \log p_\theta(x_t|c), \quad (7)$$

$$\epsilon_\theta(x_t, t, \emptyset) \propto -\nabla_{x_t} \log p_\theta(x_t), \quad (8)$$

where is an additional empty class introduced in common practices. During training, the model switches between the two modes with a ratio  $\lambda$ .

For inference, the model combines the conditional and unconditional scores and guides the denoising process as

$$\tilde{\epsilon}_\theta(x_t, t, c) = \epsilon_\theta(x_t, t, c) + w \cdot (\epsilon_\theta(x_t, t, c) - \epsilon_\theta(x_t, t, \emptyset)), \quad (9)$$

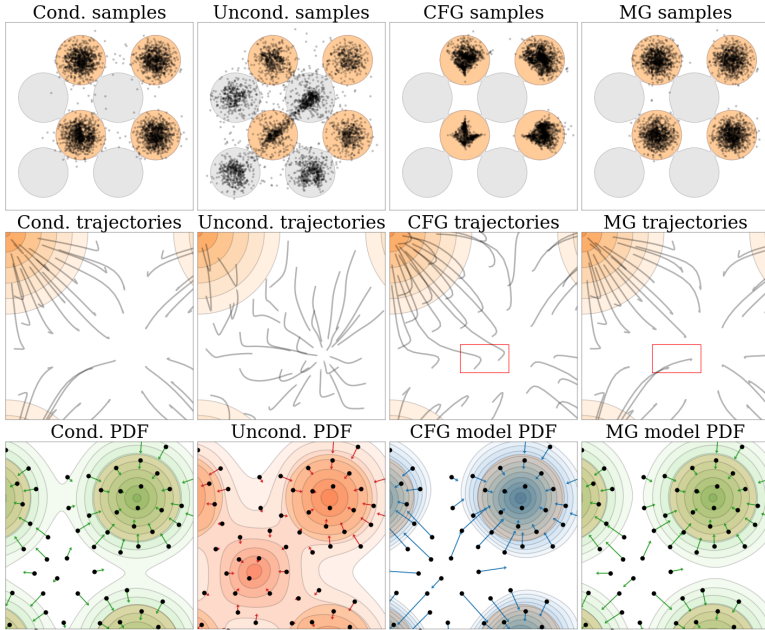


Figure 2: We use a grid 2D distribution with two classes, marked with orange and gray regions, as example and train diffusion models on it. We plot the generated samples, trajectories, and probability density function (PDF) of conditional, unconditional, CFG-guided model, and our approach.

(a) The first row indicates that although CFG improves quality by eliminating outliers, the samples concentrate in the center of data distributions, resulting the loss of diversity. In contrast, our method yields less outliers than the conditional model and a better coverage of data than CFG.

(b) In the second row, the trajectories of CFG show sharp turns at the beginning, *e.g.* samples inside the red box, while our method directly drives the samples to the closet data distributions.

(c) The PDF plots of the last row also suggest that our method predicts more symmetric contours than CFG, balancing both quality and diversity.

where  $w$  is the guidance scale that controls the focus on conditional scores and the trade-off between generation performance and sampling diversity. CFG has become an widely adopted protocol in most of diffusion models for tasks, such as image generation and video generation.

### 2.3. Distillation-based Methods

Besides acceleration (Song et al., 2023), researchers (Sauer et al., 2024) also adopt distillation on diffusion models with CFG to improve sampling quality. Rectified Flow (Liu et al., 2023) disentangles generation trajectories and streamline learning difficulty by alternatively using offline model to provide training pairs for online models. Distillation is also used to learn a smaller one-step model to match the generation performance of larger multi-step models (Meng et al., 2023). Pioneering diffusion models (Black-Forest-Labs, 2024; Stability-AI, 2024) are released with a distilled version, where CFG scale is viewed as an additional embedding to provide accurate control. However, these approaches involve two-stage learning and require extra computation and storage for offline teacher models.

## 3. Method

### 3.1. Rethinking Classifier-free guidance

Due to the complex nature of visual datasets, diffusion models often struggle whether to recover real image distribution or engage in the alignment to conditions. Classifier-free guidance (CFG) is then proposed and has become an indispensable ingredient of modern diffusion models (Nichol & Dhariwal, 2021; Karras et al., 2022; Saharia et al., 2022;

Hoogeboom et al., 2023). It drives the sample towards the regions with higher likelihood of conditions with Equation (9), where the images are more canonical and better modeled by networks (Karras et al., 2024).

However, CFG has with several disadvantages (Karras et al., 2024; Kynkäänniemi et al., 2024), such as the multitask learning of both conditional and unconditional generation, and the doubled number of function evaluations (NFEs) during inference. Moreover, the tempting property that solving the denoising process according to Equation (9) eventually recovers data distribution does not hold, as the joint distribution does not represent a valid heat diffusion of the ground-truth (Zheng & Lan, 2024). This results in exaggerated truncation and mode dropping similar to (Karras et al., 2018; Brock et al., 2019; Sauer et al., 2022), since the samples are blindly pushed towards the regions with higher posterior probability. The generation trajectories are distorted in Section 1, the images are often over-saturated in color, and the content of samples is overly simplified.

CFG originates from the classifier-guidance (Dhariwal & Nichol, 2021) that incorporates an auxiliary classifier model  $p_{\theta}(c|x_t)$  to modify the sampling distribution as

$$\tilde{p}_{\theta}(x_t|c) \propto p_{\theta}(x_t|c)p_{\theta}(c|x_t)^w, \quad (10)$$

and estimates the posterior probability term with Bayes' rule

$$p_{\theta}(c|x_t) = \frac{p_{\theta}(x_t|c)p_{\theta}(c)}{p_{\theta}(x_t)}, \quad (11)$$

where  $p_{\theta}(x_t|c)$  and  $p_{\theta}(x_t)$  are conditional and unconditional distributions, respectively.

The unconditional model is usually implemented by ran-

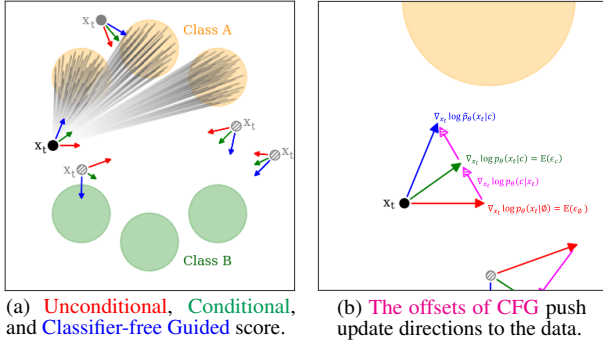


Figure 3: Illustration of our method. (a) The green and red arrow point towards the centroids of data distributions, as the training pairs  $(x_0, \epsilon)$  are randomly sampled. (b) While CFG provides accurate directions by subtracting the two vectors, our method directly learns the blue arrow,  $\nabla \log \tilde{p}_\theta(x_t|c)$ .

domly replacing labels by an empty class with a ratio  $\lambda$ . During inference, each sample is typically forwarded twice, one with and one without conditions. The finding naturally leads us to the question: can we fuse the auxiliary classifier into diffusion models in a more *efficient* and *elegant* way?

### 3.2. Model-guidance Loss

Conditional diffusion models optimize the conditional probability  $p_\theta(x_t|c)$  by Equation (3), where  $x_t$  is the noisy data and  $c$  is the condition, *e.g.*, labels and prompts. However, the models tend to ignore the condition in common practices and CFG (Ho et al., 2020) is proposed as an explicit bias.

To enhance both generation quality and alignment to conditions, we propose to take into account the posterior probability  $p_\theta(c|x_t)$ . This leads to the joint optimization of  $\tilde{p}_\theta(x_t|c) = p_\theta(x_t|c)p_\theta(c|x_t)^w$ , where  $w$  is the weighting factor of posterior probability. The score of the joint distribution is formulated as

$$\nabla_{x_t} \log \tilde{p}_\theta(x_t|c) = \nabla_{x_t} \log p_\theta(x_t|c) + w \cdot \nabla_{x_t} \log p_\theta(c|x_t) \quad (12)$$

The first term corresponds to the standard diffusion objective in Equation (3). However, the second term represents the score of posterior probability  $p_\theta(c|x_t)$  and cannot be directly obtained, since an explicit classifier of noisy samples is unavailable. Inspired by Equation (11), we transform the diffusion model into an implicit classifier and let it guide itself. Specifically, we employ Bayes' rule to estimate

$$\begin{aligned} \log p_\theta(c|x_t) &= \log p_\theta(x_t|c) - \log p_\theta(x_t) + \log p_\theta(c) \\ &\propto \log p_\theta(x_t|c) - \log p_\theta(x_t) \end{aligned} \quad (13)$$

Next, we use the diffusion model to approximate the scores

$$\nabla_{x_t} \log p_t(x_t|c) = -\frac{1}{\sigma_t} \epsilon_\theta(x_t, t, c), \quad (14)$$

$$\nabla_{x_t} \log p_t(x_t) = -\frac{1}{\sigma_t} \epsilon_\theta(x_t, t, \emptyset), \quad (15)$$

### Algorithm 1 Training with Model-guidance Loss

---

**Input:** dataset  $\{\mathbf{X}_i, \mathbf{C}_i\}$ , noise schedule  $\bar{\alpha}$ , model  $\epsilon_\theta$   
**repeat**  
 Sample data  $(x_0, c) \sim \{\mathbf{X}_i, \mathbf{C}_i\}$   
 Sample noise  $\epsilon \sim \mathcal{N}(0, 1)$  and time  $t \sim \mathbf{U}(0, 1)$   
 Add noise with  $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$   
**Modify target**  $\epsilon' = \epsilon + w \cdot \text{sg}(\epsilon_\theta(x_t, c, t) - \epsilon_\theta(x_t, \emptyset, t))$   
 Compute loss  $\mathcal{L}_{\text{MG}} = \|\epsilon_\theta(x_t, c, t) - \epsilon'\|^2$   
 Back propagation  $\theta = \theta - \eta \nabla_\theta \mathcal{L}_{\text{MG}}$   
**until** converged

---

where  $\sigma_t$  is the variance of the noise added to  $x_t$  at timestep  $t$ ,  $\emptyset$  is the empty class, and  $\epsilon_\theta(\cdot)$  is the diffusion model. Substituting Equations (14) and (15) into Equation (13) yields the score of posterior probability

$$\nabla_{x_t} \log p_\theta(c|x_t) \propto \frac{1}{\sigma_t} (\epsilon_\theta(x_t, t, \emptyset) - \epsilon_\theta(x_t, t, c)). \quad (16)$$

Then, our method applies the Bayes' estimation in Equation (13) online and trains a conditional diffusion model to directly predict the score in Equation (12), instead of separately learning Equations (14) and (15) in the form of CFG. A straight-forward implementation is to adopt the objective in Equation (3) with a modified optimization target

$$\mathcal{L}_{\text{MG}} = \mathbb{E}_{t, (x_0, c), \epsilon} \|\epsilon_\theta(x_t, t, c) - \epsilon'\|^2, \quad (17)$$

$$\epsilon' = \epsilon + w \cdot \text{sg}(\tilde{\epsilon}_\theta(x_t, t, c) - \tilde{\epsilon}_\theta(x_t, t, \emptyset)). \quad (18)$$

We apply the stop gradient operation,  $\text{sg}(\cdot)$ , which is a common practice of avoiding model collapse (Grill et al., 2020). We also use the Exponential Mean Average (EMA) counterpart of the online model,  $\tilde{\epsilon}_\theta(\cdot)$ , to stabilize the training process and provide accurate estimations. For flow-based models, we have the similar objective

$$\mathcal{L}_{\text{MG}} = \mathbb{E}_{t, (x_0, c), \epsilon} \|u_\theta(x_t, t, c) - u'\|^2, \quad (19)$$

$$u' = u + w \cdot \text{sg}(u_\theta(x_t, t, c) - u_\theta(x_t, t, \emptyset)). \quad (20)$$

where  $u$  is the ground-truth flow in Equation (6).

During training, we randomly drop the condition  $c$  in Equations (17) and (19) to  $\emptyset$  with a ratio of  $\lambda$ . These formulations transform the model itself into an implicit classifier and adjust the standard training objective of diffusion model in a self-supervised manner, allowing the joint optimization of generation quality and condition alignment with the minimum modification of existing pipelines.

### 3.3. Implementation Details

With the MG formulation in Equations (17) and (19), we have adequate options in the detailed implementations, such as incorporating an additional input of the guidance scale  $w$  into networks, replacing the usage of empty class with the law of total probability, and whether to manual or automati-



cally adjust the hyper-parameters.

**Scale-aware networks.** Similar to other distillation-based methods (Frans et al., 2024), the guidance scale  $w$  can be fed into the network as an additional condition. When augmented with  $w$ -input, our models offer flexible choices of the balance between image quality and sample diversity during inference time. Note that our models require only one forward per step for all values of  $w$ , while standard CFG needs two forwards, e.g., one with condition and one without condition. In particular, we sample guidance scale from an specified interval, and the loss function are modified into the following form

$$\mathcal{L}_{\text{MG}} = \mathbb{E}_{t, (x_0, c), \epsilon, w} \|\epsilon_{\theta}(x_t, t, c, w) - \epsilon'\|^2, \quad (21)$$

$$\epsilon' = \epsilon + w \cdot \text{sg}(\epsilon_{\theta}(x_t, t, c, 1) - \epsilon_{\theta}(x_t, t, \emptyset, 0)). \quad (22)$$

**Removing the empty class.** Another option is whether to perform multitask learning of both conditional and unconditional generation with the same model. In CFG, the estimator in Equation (11) requires to train an unconditional model. However, the multitask learning can distract and hinder model capability. Using the law of total probability

$$\begin{aligned} \nabla_{x_t} \log p_t(x_t) &= \nabla_{x_t} \log \sum_c p_t(x_t|c)p_t(c) \\ &= -\frac{1}{N\sigma_t} \sum_{i=1}^N \epsilon_{\theta}(x_t, t, c_i), \end{aligned} \quad (23)$$

where  $N$  different labels are used to estimate the unconditional score, our models focus on the conditional prediction and avoid the introduction of additional empty class.

**Automatic adjustment of the hyper-parameter  $w$ .** While the scale  $w$  in Equations (18) and (20) plays an important role, it is tedious and costly to perform manual search during training. Therefore, we introduce an automatic scheme to adjust  $w$ . We begin with  $w = 0$  that corresponds to vanilla diffusion models, then update the value with EMA according to intermediate evaluation results. The value of  $w$  is raised when quality decreases and suppressed otherwise, leading to an optimum when the training converged.

## 4. Experiment

We first present a system-level comparison with state-of-the-art models on ImageNet  $256 \times 256$  conditional generation. Then we conduct ablation experiments to investigate the detailed designs of our method. Especially, we emphasize on the following questions:

- How far can MG push the performances of existing diffusion models? (Tables 1 and 2, Section 4.2)
- How does implementation details influence the gain of proposed method? (Tables 3 to 6, Section 4.3)
- Can MG scales to larger models and datasets with efficiency? (Tables 7 and 8, Figures 4 to 6, Section 4.3)

Table 1: Experiments on ImageNet 256 without CFG. By deploying our method, the performances of both DiT-XL/2 and SiT-XL/2 are greatly boosted, achieving state-of-the-art.

MODEL	FID↓	sFID↓	IS↑	PRE.↑	REC.↑	IMG/S↑
ADM	10.9	-	101.0	0.69	0.63	-
VDM++	2.40	-	225.3	0.78	0.66	-
LDM-4	10.5	-	103.5	0.71	0.62	-
U-ViT-H	8.97	-	136.7	0.69	0.63	-
MDTV2	5.06	-	155.6	0.72	0.66	0.2
REPA	5.90	6.33	162.1	0.71	0.56	0.76
L-DiT	2.17	4.36	205.6	0.77	0.65	0.06
VAR <sub>d30</sub>	2.16	-	288.7	0.81	0.61	11.2
RAR <sub>XXL</sub>	3.83	-	274.5	0.79	0.61	3.9
MAR-H	2.35	-	227.8	0.79	0.62	0.6
DiT-XL/2	9.62	6.85	121.5	0.67	0.67	0.2
+MG <sub>(ours)</sub>	<b>2.03</b>	4.36	292.1	0.81	0.66	0.2
IMPROVE	78.9%	36.4%	140%	20.9%	1.49%	0.0%
SiT-XL/2	8.61	6.32	131.7	0.68	0.67	0.76
+MG <sub>(ours)</sub>	<b>1.34</b>	4.58	321.5	0.81	0.65	0.76
IMPROVE	84.4%	27.5%	144%	19.1%	2.99%	0.0%

Table 2: Experiments on ImageNet 256 with CFG. Comparing to models with CFG, our method still obtains excellent results and surpasses others without efficiency loss.

MODEL	FID↓	sFID↓	IS↑	PRE.↑	REC.↑	IMG/S↑
ADM	4.59	5.25	186.7	0.82	0.52	-
VDM++	2.12	-	267.7	0.81	0.65	-
LDM	3.60	-	247.7	0.87	0.48	-
U-ViT-H	2.29	5.68	263.9	0.82	0.57	-
MDTV2	1.58	4.52	314.7	0.79	0.65	0.1
REPA	1.42	4.70	305.7	0.80	0.65	0.39
L-DiT	1.35	4.15	295.3	0.79	0.65	0.03
VAR <sub>d30</sub>	1.73	-	350.2	0.82	0.60	6.3
RAR <sub>XXL</sub>	1.48	-	326.0	0.80	0.63	2.1
MAR-H	1.55	-	303.7	0.81	0.62	0.3
DiT-XL/2	2.27	4.60	278.2	0.83	0.57	0.1
+MG <sub>(ours)</sub>	<b>2.03</b>	4.36	292.1	0.81	0.66	0.2
IMPROVE	10.6%	5.22%	5.00%	2.41%	15.8%	100%
SiT-XL/2	2.06	4.49	277.5	0.83	0.59	0.39
+MG <sub>(ours)</sub>	<b>1.34</b>	4.58	321.5	0.81	0.65	0.76
IMPROVE	35.0%	2.00%	15.9%	2.41%	10.2%	94.9%

### 4.1. Setup

**Implementation and dataset.** We follow the experiment pipelines in DiT (Peebles & Xie, 2023) and SiT (Ma et al., 2024). We use ImageNet (Deng et al., 2009; Russakovsky et al., 2015) dataset and the Stable Diffusion (Rombach et al., 2022) VAE to encode  $256 \times 256$  images into the latent space of  $\mathbb{R}^{32 \times 32 \times 4}$ . We conduct ablation experiments with the B/2 variant of DiT and SiT models and train for 400K iterations. During training, we use AdamW (Kingma, 2014; Loshchilov, 2019) optimizer and a batch size of 256 in consistent with DiT (Peebles & Xie, 2023) and SiT (Ma et al., 2024) for fair comparisons. For inference, we use 1000 sampling steps for DiT models and Euler-Maruyama sampler with 250 steps for SiT.

**Baseline Models.** We compare with several state-of-the-art image generation models, including both diffusion-based and AR-based methods, which can be classified into the following three classes: (a) *Pixel-space diffusion*:

Table 3: Experiments on scale  $w$ .

MODEL	$w$	FID $\downarrow$	sFID $\downarrow$	IS $\uparrow$	PRE. $\uparrow$	REC. $\uparrow$
DiT-B/2	1.00	43.5	36.7	39.23	0.62	0.34
+MG <sub>(ours)</sub>	1.25	9.86	8.87	176.1	0.81	0.37
+MG <sub>(ours)</sub>	1.50	<b>7.24</b>	<b>5.56</b>	189.2	0.84	0.38
+MG <sub>(ours)</sub>	1.75	8.21	6.63	197.2	<b>0.86</b>	0.38
+MG <sub>(ours)</sub>	2.00	9.66	7.90	<b>224.7</b>	0.85	<b>0.39</b>
+MG <sub>(ours)</sub>	AUTO	7.60	6.29	192.4	0.85	0.38
SiT-B/2	1.00	33.0	27.8	65.24	0.68	0.35
+MG <sub>(ours)</sub>	1.25	8.94	7.87	194.3	0.83	0.38
+MG <sub>(ours)</sub>	1.50	<b>6.49</b>	<b>5.69</b>	212.3	0.86	0.38
+MG <sub>(ours)</sub>	1.75	8.03	6.91	221.0	0.86	0.39
+MG <sub>(ours)</sub>	2.00	9.14	7.99	<b>236.7</b>	<b>0.88</b>	<b>0.40</b>
+MG <sub>(ours)</sub>	AUTO	6.86	5.88	219.1	0.87	0.38

 Table 4: Experiments on drop ratio  $\lambda$ .

MODEL	$\lambda$	FID $\downarrow$	sFID $\downarrow$	IS $\uparrow$	PRE. $\uparrow$	REC. $\uparrow$
DiT-B/2	1.00	43.5	36.7	39.23	0.62	0.34
+MG <sub>(ours)</sub>	0.05	11.7	9.90	156.7	0.78	0.33
+MG <sub>(ours)</sub>	0.10	<b>7.24</b>	<b>5.56</b>	<b>189.2</b>	<b>0.84</b>	<b>0.38</b>
+MG <sub>(ours)</sub>	0.15	7.62	5.99	183.4	0.83	0.38
+MG <sub>(ours)</sub>	0.20	9.01	7.04	171.7	0.81	0.36
SiT-B/2	1.00	33.0	27.8	65.24	0.68	0.35
+MG <sub>(ours)</sub>	0.05	10.8	9.25	168.8	0.80	0.34
+MG <sub>(ours)</sub>	0.10	<b>6.49</b>	<b>5.69</b>	<b>212.3</b>	<b>0.86</b>	<b>0.38</b>
+MG <sub>(ours)</sub>	0.15	6.77	5.89	207.4	0.85	0.37
+MG <sub>(ours)</sub>	0.20	8.87	8.06	199.6	0.84	0.37

ADM (Dhariwal & Nichol, 2021), VDM++ (Kingma & Gao, 2023); (b) *Latent-space diffusion*: LDM (Rombach et al., 2022), U-ViT (Bao et al., 2023), MDTv2 (Gao et al., 2023), REPA (Yu et al., 2024b), LightningDiT(L-DiT) (Yao & Wang, 2025), DiT (Peebles & Xie, 2023), SiT (Ma et al., 2024); (c) *Auto-regressive models*: VAR (Tian et al., 2024), RAR (Yu et al., 2024a), MAR (Li et al., 2024). These models consist of strong baselines and demonstrate the advantages of our method. Although our method does not require CFG during inference, we still compare with these baselines under two settings, with and without CFG, for thoroughly investigations.

**Evaluation metrics.** We report the commonly used Frechet inception distance (Heusel et al., 2017) with 50,000 samples (FID-50K). In addition, we report sFID (Nash et al., 2021), Inception Score (IS) (Salimans et al., 2016), Precision (Pre.), and Recall (Rec.) (Kynkäänniemi et al., 2019) as supplementary metrics. We also report the time to generate one sample of each model in seconds to measure the trade-off between generation quality and computation budget.

## 4.2. Overall Performances

First of all, we present a through system-level comparison with recent state-of-the-art image generation approaches on ImageNet  $256 \times 256$  dataset in Tables 1 and 2. As shown in Table 1, both DiT-XL/2 and SiT-XL/2 models greatly benefit from our method, achieving the outstanding performance gain of 78.9% and 84.4%. It is worth mentioning that our models do not apply modern techniques in the inference process, including rejection sampling (Tian et al.,

 Table 5: Experiments on Model input  $w$ .

MODEL	$w$ -IN	FID $\downarrow$	sFID $\downarrow$	IS $\uparrow$	PRE. $\uparrow$	REC. $\uparrow$
DiT-B/2	$\times$	43.5	36.7	39.23	0.62	0.34
+MG <sub>(ours)</sub>	$\times$	<b>7.24</b>	<b>5.56</b>	<b>189.2</b>	<b>0.84</b>	0.38
+MG <sub>(ours)</sub>	$\checkmark$	8.13	6.03	175.1	0.84	<b>0.39</b>
SiT-B/2	$\times$	33.0	27.8	65.24	0.68	0.35
+MG <sub>(ours)</sub>	$\times$	<b>6.49</b>	<b>5.69</b>	<b>212.3</b>	<b>0.86</b>	<b>0.38</b>
+MG <sub>(ours)</sub>	$\checkmark$	7.33	5.96	207.4	0.85	0.38

 Table 6: Experiments on empty class  $\emptyset$ .

MODEL	$\emptyset$ -CLS	FID $\downarrow$	sFID $\downarrow$	IS $\uparrow$	PRE. $\uparrow$	REC. $\uparrow$
DiT-B/2	$\times$	43.5	36.7	39.23	0.62	0.34
+MG <sub>(ours)</sub>	$\times$	9.66	8.73	174.4	0.81	0.35
+MG <sub>(ours)</sub>	$\checkmark$	<b>7.24</b>	<b>5.56</b>	<b>189.2</b>	<b>0.84</b>	<b>0.38</b>
SiT-B/2	$\times$	33.0	27.8	65.24	0.68	0.35
+MG <sub>(ours)</sub>	$\times$	9.03	7.96	183.3	0.82	0.35
+MG <sub>(ours)</sub>	$\checkmark$	<b>6.49</b>	<b>5.69</b>	<b>212.3</b>	<b>0.86</b>	<b>0.38</b>

2024), classifier-free guidance (Ho et al., 2020) and guidance interval (Kynkäänniemi et al., 2024). Compared to advanced methods, our models are light-weight, *e.g.* 675M in contrast to RAR-XXL with 1.5B and MAR-H with 943M parameters, and consume less computational resources, for example, LightningDiT uses DiT-XL/1 to reduce patch size to  $1 \times 1$  and needs  $16 \times$  computation in attention operations.

To facilitate a fair evaluation, we also compare with other methods with Classifier-free guidance. While prevalent diffusion models significantly benefit and are indispensable from CFG, it introduces an additional forward without condition and doubles the computation consumptions. Also, it usually requires a careful search over the hyper-parameter of guidance scale to achieve the best trade-off between quality and diversity. In contrast, our models still surpass other CFG-assisted methods and run with only half of the generation time.

Finally, we report the time consumption for each model to generate one sample in seconds. Comparing to other diffusion-based approaches facilitated with vanilla CFG, our method runs significantly faster and does not sacrifice inference speed for sampling quality.

## 4.3. Ablation study

To thoroughly understand the designs and subsequent influences of our method, we conduct ablation experiments on the key components, including the hyper-parameter  $w$ ,  $\lambda$  choices, whether the model takes  $w$  as input, and the role of empty class during training. Moreover, we assess the scalability of our method in terms of both model size and dataset difficulty.

**Hyper-parameter  $w$**  In Equations (18) and (20), the hyper-parameter  $w$  controls the scale of posterior probability and serves an important role akin to the guidance scale in CFG, which is sensitive to FID-score. We conduct ablation experiments on the hyper-parameter  $w$  and report results in

Table 7: Experiments on Model size. Our method scales to models with different sizes.

MODEL	FID↓	sFID↓	IS↑	PRE.↑	REC.↑
DiT-B/2	43.5	36.7	39.23	0.62	0.34
+MG <sub>(ours)</sub>	7.24	5.56	189.2	0.84	0.38
DiT-L/2	23.3	18.4	132.7	0.73	0.40
+MG <sub>(ours)</sub>	5.43	4.66	236.3	0.83	0.44
DiT-XL/2	19.5	15.6	163.5	0.79	0.46
+MG <sub>(ours)</sub>	3.37	4.73	257.2	0.84	0.51
SiT-B/2	33.0	27.8	65.24	0.68	0.35
+MG <sub>(ours)</sub>	6.49	5.69	212.3	0.86	0.38
SiT-L/2	18.9	16.3	173.2	0.71	0.42
+MG <sub>(ours)</sub>	4.50	4.03	243.9	0.85	0.46
SiT-XL/2	17.3	13.9	192.1	0.78	0.50
+MG <sub>(ours)</sub>	2.89	3.12	261.0	0.85	0.54

Table 8: Experiments on ImageNet 512. Our method scales to high-resolution image datasets.

MODEL	FID↓	sFID↓	IS↑	PRE.↑	REC.↑
DiT-XL/2	3.04	5.02	240.8	0.84	0.54
+MG <sub>(ours)</sub>	2.78	4.86	257.2	0.83	0.58
SiT-XL/2	2.62	4.18	252.2	0.84	0.57
+MG <sub>(ours)</sub>	2.24	4.03	276.9	0.86	0.60

Table 3, where  $w = 1$  refers to vanilla forms in DiT (Peebles & Xie, 2023) and SiT (Ma et al., 2024). It is shown that the choice of  $w$  also acts as a crucial role and balances the trade-off between quality and diversity.

To overcome the tiresome and costly search of  $w$  during training, we propose an adaptive approach to automatically adjust  $w$ , which achieves comparable performance with manual search. Meanwhile, we can further apply CFG to our models in Figure 4 to flexibly adjust between better quality and diversity during inference.

**Hyper-parameter  $\lambda$**  The relative ratio to train conditional and unconditional models,  $\lambda$ , is also important to our method. The unconditional model is usually trained by randomly dropping the condition and replacing with an additional empty label for part of training data. In Table 4, we conduct ablation experiments on the hyper-parameter  $\lambda$  report the corresponding results. We find that  $\lambda$  is less sensitive than  $w$ , and  $\lambda \in \{0.10, 0.15\}$  offers satisfactory performances.

**Model input  $w$**  Despite the same loss formulation in the Equation (17), it is optional whether our model takes the scale  $w$  as an additional input. In Table 5, the models with  $w$ -input slightly lag behind the counterparts without  $w$ -input but still exceeding the vanilla DiT-B/2 and SiT-B/2 with CFG, demonstrating the superiority of our method.

**Empty class  $\emptyset$**  In Table 6, we conduct ablation experiments on the introduction of additional empty class. While removing the empty class in our method leads to worse estimation of posterior probability, the generation performances are still on par with the vanilla CFG. It can also be improved by

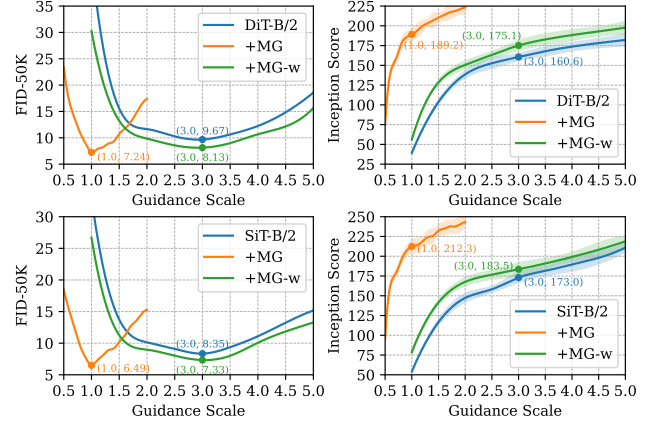


Figure 4: FID-50K and Inception Score results as the guidance scale increases during inference. Our method is compatible with and can be wrapped into vanilla CFG.

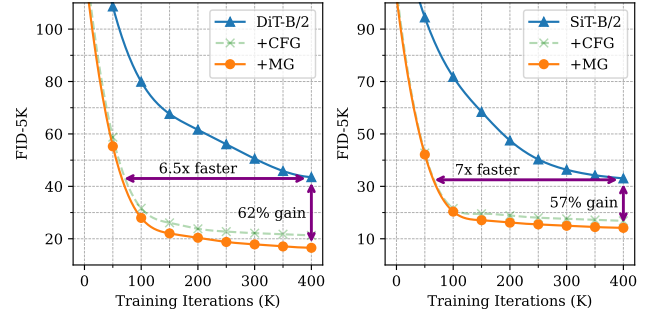


Figure 5: FID-5K results during training. Our method is  $\geq 6.5\times$  faster and  $\approx 60\%$  better than vanilla DiT and SiT, even surpassing the results of CFG.

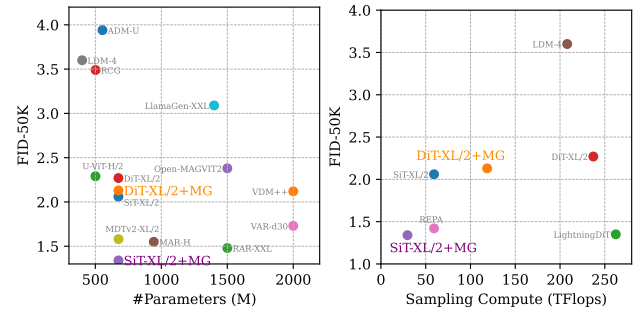
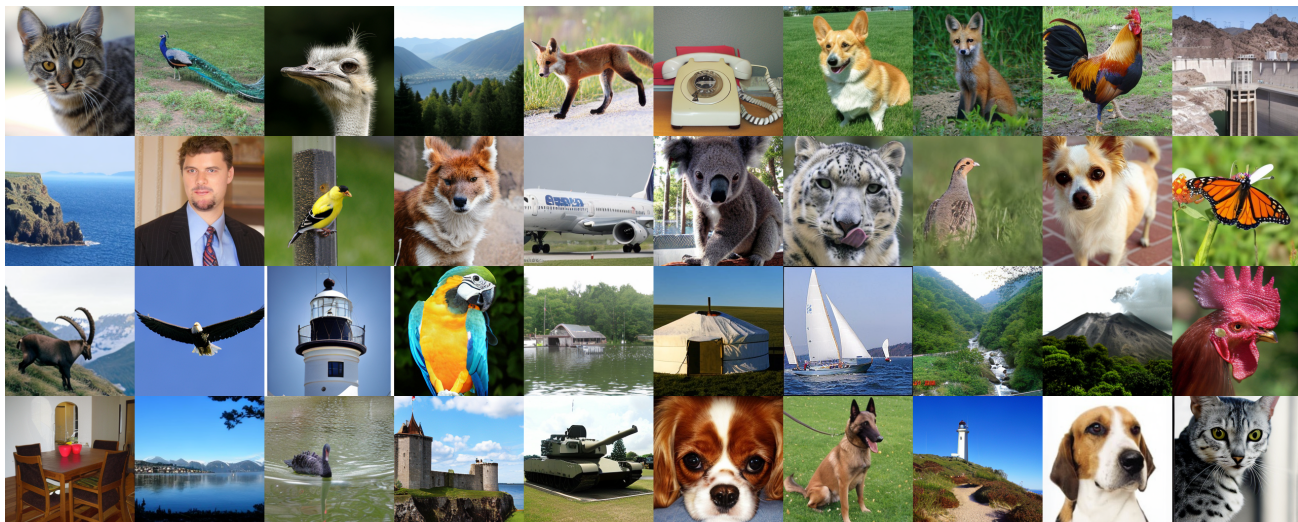


Figure 6: FID-50K vs. number of parameters and sampling flops of different models, where our models are highlighted.

better estimation with the law of total probability or a larger batch size.

**Efficiency** One key advantage of our method is that it not only improves inference speed by avoiding the second network forward of CFG, but also accelerates the training and convergence of diffusion models. In Figure 5, our method obtains  $\geq 6.5\times$  convergence speed and  $\approx 60\%$  performance gain. In Figure 6, we plot the number of network parameters and sampling compute in TFlops versus FID-50K of



Figure 7: **Uncurated** samples of SiT-XL/2+MG on ImageNet  $256 \times 256$ .Figure 8: **Uncurated** samples of SiT-XL/2+MG on ImageNet  $512 \times 512$ .

concurrent methods. When comparing number of network parameters, our method comes with the lowest FID and a small model size. When comparing sampling computes, our method achieves state-of-the-art performances in parallel with LightningDiT (Yao & Wang, 2025), while requires only  $\approx 12\%$  computational resources.

**Scalability** Finally, the scalability to larger model and dataset of our method is of imparible significance. In Table 7, we conduct ablation stuides on model size with B/2, L/2 and XL/2 variants of DiT and SiT models. It is demonstrated that our method is capable to boost the performance of models with different sizes and designs. We scale to ImageNet  $512 \times 512$  dataset to validate our method in handling difficult distributions in Table 8. As depicted, our method also offers improvements on high-resolution tasks.

## 5. Conclusion

This work addresses the limitations of the commonly used Classifier-free guidance (CFG) of diffusion models, and proposes Model-guidance (MG) as an efficient and advantageous replacement. We first investigate the mechanism of CFG and locate the source of performance gain as a joint optimization of posterior probability. Then, we transcend the idea into the training process of diffusion models and directly learn the score of the joint distribution,  $\nabla \log \tilde{p}_\theta(x_t|c) = \nabla \log p_\theta(x_t|c)p_\theta(c|x_t)^w$ . Comprehensive experiments demonstrate that our method significantly boosts the generation performance without efficiency loss, scales to different models and datasets, and achieves state-of-the-art results on ImageNet  $256 \times 256$  dataset. We believe that this work contributes to future diffusion models.



## Impact Statements

This paper propose methods in association with generative methods. There might be potential negative social impacts, *e.g.* generating fake portraits, as the core contribution of our work is a new algorithm of generative modeling. As possible mitigation strategies, we will restrict the access to these models in the planned release of code and models. We also validate that current detectors can effectively determine our generation results about human portraits.

## References

- Albergo, M. S., Boffi, N. M., and Vanden-Eijnden, E. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- Bao, F., Nie, S., Xue, K., Cao, Y., Li, C., Su, H., and Zhu, J. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22669–22679, 2023.
- Black-Forest-Labs. Flux.1 model family, 2024. URL <https://blackforestlabs.ai/announcing-black-forest-labs/>.
- Blattmann, A., Rombach, R., Ling, H., Dockhorn, T., Kim, S. W., Fidler, S., and Kreis, K. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22563–22575, 2023.
- Brock, A., Donahue, J., and Simonyan, K. Large scale gan training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2019.
- Chen, J., Jincheng, Y., Chongjian, G., Yao, L., Xie, E., Wang, Z., Kwok, J., Luo, P., Lu, H., and Li, Z. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Dhariwal, P. and Nichol, A. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Frans, K., Hafner, D., Levine, S., and Abbeel, P. One step diffusion via shortcut models. *arXiv preprint arXiv:2410.12557*, 2024.
- Gao, S., Zhou, P., Cheng, M.-M., and Yan, S. Masked diffusion transformer is a strong image synthesizer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 23164–23173, 2023.
- Grill, J.-B., Strub, F., Althché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- Gupta, A., Yu, L., Sohn, K., Gu, X., Hahn, M., Li, F.-F., Essa, I., Jiang, L., and Lezama, J. Photorealistic video generation with diffusion models. In *European Conference on Computer Vision*, pp. 393–411. Springer, 2025.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- Ho, J. and Salimans, T. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Ho, J., Salimans, T., Gritsenko, A., Chan, W., Norouzi, M., and Fleet, D. J. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.
- Hoogeboom, E., Heek, J., and Salimans, T. simple diffusion: End-to-end diffusion for high resolution images. In *International Conference on Machine Learning*, pp. 13213–13232. PMLR, 2023.
- Karras, T., Laine, S., and Aila, T. A style-based generator architecture for generative adversarial networks. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4396–4405, 2018. URL <https://api.semanticscholar.org/CorpusID:54482423>.
- Karras, T., Aittala, M., Aila, T., and Laine, S. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35: 26565–26577, 2022.
- Karras, T., Aittala, M., Kynkäänniemi, T., Lehtinen, J., Aila, T., and Laine, S. Guiding a diffusion model with a bad version of itself. *Advances in neural information processing systems*, 2024.
- Kingma, D. P. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

- Kingma, D. P. and Gao, R. Understanding the diffusion objective as a weighted integral of elbos. *arXiv preprint arXiv:2303.00848*, 2, 2023.
- Kynkäänniemi, T., Karras, T., Laine, S., Lehtinen, J., and Aila, T. Improved precision and recall metric for assessing generative models. *Advances in neural information processing systems*, 32, 2019.
- Kynkäänniemi, T., Aittala, M., Karras, T., Laine, S., Aila, T., and Lehtinen, J. Applying guidance in a limited interval improves sample and distribution quality in diffusion models. *Advances in neural information processing systems*, 2024.
- Li, T., Tian, Y., Li, H., Deng, M., and He, K. Autoregressive image generation without vector quantization. *Advances in neural information processing systems*, 2024.
- Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and Le, M. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- Liu, X., Gong, C., et al. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2023.
- Loshchilov, I. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- Ma, N., Goldstein, M., Albergo, M. S., Boffi, N. M., VandenEijnden, E., and Xie, S. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024.
- McCann, R. J. A convexity principle for interacting gases. *Advances in mathematics*, 128(1):153–179, 1997.
- Meng, C., Rombach, R., Gao, R., Kingma, D., Ermon, S., Ho, J., and Salimans, T. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14297–14306, 2023.
- Nash, C., Menick, J., Dieleman, S., and Battaglia, P. Generating images with sparse representations. In *International Conference on Machine Learning*, pp. 7958–7968. PMLR, 2021.
- Nichol, A. Q. and Dhariwal, P. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.
- Nichol, A. Q., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., Mcgrew, B., Sutskever, I., and Chen, M. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, pp. 16784–16804. PMLR, 2022.
- Peebles, W. and Xie, S. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., and Rombach, R. Sdxl: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024.
- Polyak, A., Zohar, A., Brown, A., Tjandra, A., Sinha, A., Lee, A., Vyas, A., Shi, B., Ma, C.-Y., Chuang, C.-Y., et al. Movie gen: A cast of media foundation models. *arXiv preprint arXiv:2410.13720*, 2024.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115: 211–252, 2015.
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E. L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35: 36479–36494, 2022.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016.
- Sauer, A., Schwarz, K., and Geiger, A. Stylegan-xl: Scaling stylegan to large diverse datasets. *ACM SIGGRAPH 2022 Conference Proceedings*, 2022. URL <https://api.semanticscholar.org/CorpusID:246441861>.
- Sauer, A., Boesel, F., Dockhorn, T., Blattmann, A., Esser, P., and Rombach, R. Fast high-resolution image synthesis with latent adversarial diffusion distillation. In *SIGGRAPH Asia 2024 Conference Papers*, pp. 1–11, 2024.

- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pp. 2256–2265. PMLR, 2015.
- Song, J., Meng, C., and Ermon, S. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021a.
- Song, Y. and Ermon, S. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021b.
- Song, Y., Dhariwal, P., Chen, M., and Sutskever, I. Consistency models. In *International Conference on Machine Learning*, pp. 32211–32252. PMLR, 2023.
- Stability-AI. Introducing stable diffusion 3.5, 2024. URL <https://stability.ai/news/introducing-stable-diffusion-3-5>.
- Tian, K., Jiang, Y., Yuan, Z., Peng, B., and Wang, L. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *Advances in neural information processing systems*, 2024.
- Tong, A., FATRAS, K., Malkin, N., Huguet, G., Zhang, Y., Rector-Brooks, J., Wolf, G., and Bengio, Y. Improving and generalizing flow-based generative models with minibatch optimal transport. *Transactions on Machine Learning Research*, 2024.
- Wang, Y., Chen, X., Ma, X., Zhou, S., Huang, Z., Wang, Y., Yang, C., He, Y., Yu, J., Yang, P., et al. Lavie: High-quality video generation with cascaded latent diffusion models. *International Journal of Computer Vision*, pp. 1–20, 2024.
- Yao, J. and Wang, X. Reconstruction vs. generation: Taming optimization dilemma in latent diffusion models. *arXiv preprint arXiv:2501.01423*, 2025.
- Yu, Q., He, J., Deng, X., Shen, X., and Chen, L.-C. Randomized autoregressive visual generation. *arXiv preprint arXiv:2411.00776*, 2024a.
- Yu, S., Kwak, S., Jang, H., Jeong, J., Huang, J., Shin, J., and Xie, S. Representation alignment for generation: Training diffusion transformers is easier than you think. *arXiv preprint arXiv:2410.06940*, 2024b.
- Zheng, C. and Lan, Y. Characteristic guidance: Non-linear correction for diffusion model at large guidance scale. In *International Conference on Machine Learning*. PMLR, 2024.