

# REVAMPING DIFFUSION GUIDANCE FOR CONDITIONAL AND UNCONDITIONAL GENERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Classifier-free guidance (CFG) has become the standard method for enhancing the quality of conditional diffusion models. However, employing CFG requires either training an unconditional model alongside the main diffusion model or modifying the training procedure by periodically inserting a null condition. There is also no clear extension of CFG to unconditional models. In this paper, we revisit the core principles of CFG and introduce a new method, independent condition guidance (ICG), which provides the benefits of CFG without the need for any special training procedures. Our approach streamlines the training process of conditional diffusion models and can also be applied during inference on any pre-trained conditional model. Additionally, by leveraging the time-step information encoded in all diffusion networks, we propose an extension of CFG, called time-step guidance (TSG), which can be applied to *any* diffusion model, including unconditional ones. Our guidance techniques are easy to implement and have the same sampling cost as CFG. Through extensive experiments, we demonstrate that ICG matches the performance of standard CFG across various conditional diffusion models. Moreover, we show that TSG improves generation quality in a manner similar to CFG, without relying on any conditional information.

## 1 INTRODUCTION

Diffusion models have recently emerged as the main methodology behind many successful generative models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Dhariwal & Nichol, 2021; Rombach et al., 2022; Song & Ermon, 2019; Song et al., 2021b). At the core of such models lies a diffusion process that gradually adds noise to the data until the corrupted points are indistinguishable from pure noise. During inference, a denoiser trained to reverse this process is used to gradually refine pure-noise samples until they resemble the clean data. While the theory suggests that standard sampling from diffusion models should yield high-quality images, this does not generally hold in practice, and guidance methods are often required to increase the quality of generations, albeit at the expense of diversity (Dhariwal & Nichol, 2021; Ho & Salimans, 2022). Classifier guidance (Dhariwal & Nichol, 2021) introduced this quality-boosting concept by utilizing the gradient of a classifier trained on noisy images to increase the class-likelihood of generated samples. Later, classifier-free guidance (CFG) (Ho & Salimans, 2022) was proposed, allowing diffusion models to simulate the same behavior as classifier guidance without using an explicit classifier. Since then, CFG has been applied to other conditional generation tasks, such as text-to-image synthesis (Nichol et al., 2022) and text-to-3D generation (Poole et al., 2023).

In addition to CFG’s trading diversity for quality, it has the following two practical limitations. First, it requires a dedicated, pre-defined training process on an *auxiliary* task in order to learn the unconditional score function. This typically involves training a separate unconditional model or, more commonly, randomly dropping the conditioning vector and replacing it with a null vector during training. This approach reduces training efficiency, as the model now needs to be trained on two different tasks. Moreover, replacing the condition may not be straightforward when multiple conditioning signals—such as text, images, and audio—are used simultaneously or when the null vector (often the zero vector) carries specific meaning. We demonstrate that this dedicated auxiliary training process is unnecessary. A second limitation is that there has been no clear way to extend the benefits of classifier-free guidance beyond conditional models to unconditional generation. We introduce a method that closes this gap.

054 We revisit the principles behind classifier-free guidance and show both theoretically and empirically  
055 that similar quality-boosting behavior can be achieved without the need for additional auxiliary  
056 training of an unconditional model. The main idea is that by using a conditioning vector *independent*  
057 of the input data, the conditional score function becomes equivalent to the unconditional score. This  
058 insight leads us to propose *independent condition guidance* (ICG), a method that replicates the  
059 behavior of CFG at inference time without requiring auxiliary training of an unconditional model,  
060 i.e., without needing explicit access to the unconditional score function. In Section 6.1, we show that  
061 the auxiliary training of the unconditional model in CFG can be detrimental to training efficiency,  
062 and similar or better performance can be achieved by training only a purely conditional model and  
063 using ICG instead.

064 Inspired by the above, we also introduce a novel technique to extend classifier-free guidance to a more  
065 general setting that includes unconditional generation. We argue that by using a perturbed version  
066 of the time-step embedding in diffusion models, one can create a guidance signal similar to CFG to  
067 improve the quality of generations. This method, which we call *time-step guidance* (TSG), aims to  
068 improve the accuracy of denoising at each sampling step by leveraging the time-step information  
069 learned by the diffusion model to steer sampling trajectories toward better noise-removal paths.

070 ICG and TSG are easy to implement, do not require additional fine-tuning of the underlying diffusion  
071 models, and have the same sampling cost as CFG. Through extensive experiments, we empirically  
072 verify that: 1) ICG offers performance similar to CFG and can be readily applied to models that are  
073 not trained with the CFG objective in mind, such as EDM (Karras et al., 2022); and 2) TSG improves  
074 output quality in a manner similar to CFG for both conditional and unconditional generation.

075 The core contributions of our work are as follows: (i) We revisit the principles of classifier-free  
076 guidance and offer an efficient, theoretically motivated method to employ CFG without requiring any  
077 auxiliary training of an unconditional model, greatly simplifying the training process of conditional  
078 diffusion models and improving training efficiency relative to the standard approach. (ii) We offer an  
079 extension of CFG that is generally applicable to all diffusion models, whether conditional or uncondi-  
080 tional. (iii) We demonstrate empirically that our guidance techniques achieve the quality-boosting  
081 benefits of CFG across various setups and network architectures.

## 082 083 2 RELATED WORK 084

085 Score-based diffusion models (Song & Ermon, 2019; Song et al., 2021b; Sohl-Dickstein et al., 2015;  
086 Ho et al., 2020) learn the data distribution by reversing a forward diffusion process that progressively  
087 transforms the data into Gaussian noise. These models have quickly surpassed the fidelity and  
088 diversity of previous generative modeling methods (Nichol & Dhariwal, 2021; Dhariwal & Nichol,  
089 2021), achieving state-of-the-art results in various domains, including unconditional image generation  
090 (Dhariwal & Nichol, 2021; Karras et al., 2022), text-to-image generation (Ramesh et al., 2022;  
091 Saharia et al., 2022b; Balaji et al., 2022; Rombach et al., 2022; Podell et al., 2023; Yu et al., 2022),  
092 video generation (Blattmann et al., 2023b;a; Gupta et al., 2023), image-to-image translation (Saharia  
093 et al., 2022a; Liu et al., 2023), motion synthesis (Tevet et al., 2023; Tseng et al., 2023), and audio  
094 generation (Chen et al., 2021; Kong et al., 2021; Huang et al., 2023).

095 Since the development of the DDPM model (Ho et al., 2020), many advancements have been proposed  
096 including improved network architectures (Hoogeboom et al., 2023; Karras et al., 2023; Peebles &  
097 Xie, 2022; Dhariwal & Nichol, 2021), sampling algorithms (Song et al., 2021a; Karras et al., 2022;  
098 Liu et al., 2022; Lu et al., 2022a; Salimans & Ho, 2022), and training methods (Nichol & Dhariwal,  
099 2021; Karras et al., 2022; Song et al., 2021b; Salimans & Ho, 2022; Rombach et al., 2022). Despite  
100 these recent advances, diffusion guidance, including classifier and classifier-free guidance (Dhariwal  
101 & Nichol, 2021; Ho & Salimans, 2022), still plays an essential role in improving the quality of  
102 generations as well as increasing the alignment between the condition and the output image (Nichol  
103 et al., 2022).

104 SAG (Hong et al., 2022) and PAG (Ahn et al., 2024) have recently been proposed to increase the  
105 quality of UNet-based diffusion models by modifying the predictions of the self-attention layers. Our  
106 method is complementary to these approaches, as one can combine ICG updates with the update  
107 signal from the perturbed attention modules (Hong et al., 2022). In addition, we make no assumptions  
about the network architecture.

Another line of work includes guiding the generation of the diffusion model with a differentiable loss function or an off-the-shelf classifier (Song et al., 2023; Chung et al., 2022; Yu et al., 2023; Bansal et al., 2023; He et al., 2023). These methods are primarily focused on solving inverse problems, typically with unconditional models, while we are instead concerned with achieving the benefits of CFG in conditional models without any additional training requirements. With TSG, we also generalize our approach to extend CFG-like benefits to unconditional models.

Perturbing the condition vector is employed in CADs (Sadat et al., 2024) to increase the diversity of generations. CADs differs from ICG in focusing on the *conditional* branch to improve diversity, while ICG is concerned with the *unconditional* branch to simulate CFG. Since CADs is designed to enhance the diversity of CFG, it can be used alongside ICG to improve the diversity of output at high guidance scales (see Appendix B).

### 3 BACKGROUND

This section provides an overview of diffusion models. Let  $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$  be a data point,  $t \in [0, 1]$  be the time step, and  $\mathbf{z}_t = \mathbf{x} + \sigma(t)\boldsymbol{\epsilon}$  be the forward process of the diffusion model that adds noise to the data. Here  $\sigma(t)$  is the noise schedule and determines how much information is destroyed at each time step  $t$ , with  $\sigma(0) = 0$  and  $\sigma(1) = \sigma_{\text{max}}$ . Karras et al. (2022) showed that this forward process corresponds to the ordinary differential equation (ODE)

$$d\mathbf{z}_t = -\dot{\sigma}(t)\sigma(t) \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) dt \quad (1)$$

or, equivalently, a stochastic differential equation (SDE) given by

$$d\mathbf{z}_t = -\dot{\sigma}(t)\sigma(t) \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) dt - \beta(t)\sigma(t)^2 \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) dt + \sqrt{2\beta(t)}\sigma(t) d\omega_t. \quad (2)$$

Here  $d\omega_t$  is the standard Wiener process, and  $p_t(\mathbf{z}_t)$  is the time-dependent distribution of noisy samples, with  $p_0 = p_{\text{data}}$  and  $p_1 = \mathcal{N}(\mathbf{0}, \sigma_{\text{max}}^2 \mathbf{I})$ . Assuming that we have access to the time-dependent score function  $\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t)$ , we can sample from the data distribution  $p_{\text{data}}$  by solving the ODE or SDE backward in time (from  $t = 1$  to  $t = 0$ ). The unknown score function  $\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t)$  is estimated via a neural denoiser  $D_{\boldsymbol{\theta}}(\mathbf{z}_t, t)$  that is trained to predict the clean samples  $\mathbf{x}$  from the corresponding noisy samples  $\mathbf{z}_t$ . The framework allows for conditional generation by training a denoiser  $D_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{y})$  that accepts additional input signals  $\mathbf{y}$ , such as class labels or text prompts.

**Training objective** Given a noisy sample  $\mathbf{z}_t$  at time step  $t$ , the denoiser  $D_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{y})$  with parameters  $\boldsymbol{\theta}$  can be trained with the standard MSE loss (also called denoising score matching loss)

$$\arg \min_{\boldsymbol{\theta}} \mathbb{E}_t \left[ \|D_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{y}) - \mathbf{x}\|^2 \right]. \quad (3)$$

The denoiser approximates the time-dependent conditional score function  $\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t | \mathbf{y})$  via

$$\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t | \mathbf{y}) \approx \frac{D_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{y}) - \mathbf{z}_t}{\sigma(t)^2}. \quad (4)$$

**Classifier-free guidance (CFG)** CFG is an inference method for improving the quality of generated outputs by mixing the predictions of a conditional and an unconditional model (Ho & Salimans, 2022). Specifically, given a null condition  $\mathbf{y}_{\text{null}} = \emptyset$  corresponding to the unconditional case, CFG modifies the output of the denoiser at each sampling step according to

$$\hat{D}_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{y}) = D_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{y}_{\text{null}}) + w_{\text{CFG}}(D_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{y}) - D_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{y}_{\text{null}})), \quad (5)$$

where  $w_{\text{CFG}} = 1$  corresponds to the non-guided case. The unconditional model  $D_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{y}_{\text{null}})$  is trained by randomly assigning the null condition  $\mathbf{y}_{\text{null}} = \emptyset$  to the input of the denoiser with probability  $p$ , where we normally have  $p \in [0.1, 0.2]$ . One can also train a separate denoiser to estimate the unconditional score in Equation (5) (Karras et al., 2023). Similar to the truncation method in GANs (Brock et al., 2019), CFG increases the quality of individual images at the expense of less diversity (Murphy, 2023).

## 4 REVISITING CLASSIFIER-FREE GUIDANCE

We now show how a conditional model can be used to simulate the behavior of classifier-free guidance, without needing any auxiliary training to learn the unconditional score function. The analysis in this section is inspired by Sadat et al. (2024).

First, note that at each time step  $t$ , classifier-free guidance uses the conditional score  $\nabla_{z_t} \log p_t(z_t | \mathbf{y})$  and the unconditional score  $\nabla_{z_t} \log p_t(z_t)$  to guide the sampling process. Based on Bayes' theorem, we can write  $p_t(z_t | \mathbf{y}) = \frac{p_t(\mathbf{y} | z_t) p_t(z_t)}{p_t(\mathbf{y})}$ , which gives us

$$\nabla_{z_t} \log p_t(z_t | \mathbf{y}) = \nabla_{z_t} \log p_t(z_t) + \nabla_{z_t} \log p_t(\mathbf{y} | z_t). \quad (6)$$

Next, assume that we replace the condition with a random vector  $\hat{\mathbf{y}}$  that is independent of the input  $z_t$ . In this case, we have  $p_t(\hat{\mathbf{y}} | z_t) = p_t(\hat{\mathbf{y}})$ , which gives us

$$\nabla_{z_t} \log p_t(z_t | \hat{\mathbf{y}}) = \nabla_{z_t} \log p_t(z_t) + \nabla_{z_t} \log p_t(\hat{\mathbf{y}}) = \nabla_{z_t} \log p_t(z_t). \quad (7)$$

This analysis shows that we can estimate the unconditional score purely based on the conditional model by replacing the condition  $\mathbf{y}$  with an independent vector  $\hat{\mathbf{y}}$ . Based on this derivation, we argue that there is no need to train a separate model  $D_{\theta}(z_t, t, \mathbf{y}_{\text{null}})$  for applying classifier-free guidance as we can use the conditional model itself to bootstrap the score of the unconditional distribution as long as we pick an input condition that is independent of  $z_t$ . We call this method *independent condition guidance* (ICG) for the rest of the paper.

The intuition behind ICG becomes more clear in the class-conditional case. Notice that by knowing the conditional distribution  $p_t(z_t | \mathbf{y})$  for each  $\mathbf{y}$ , we also implicitly obtain the unconditional distribution through  $p_t(z_t) = \sum_{\mathbf{y}} p_t(z_t | \mathbf{y}) p(\mathbf{y})$ . The major limitation of directly applying this formula is the necessity for multiple forward passes (one for each class). ICG circumvents this issue by deriving the unconditional score efficiently from the conditional model using only a single forward pass through the network. Thus, the sampling cost of ICG is equal to that of standard CFG.

**Implementation details** Although the analysis above holds for any condition  $\hat{\mathbf{y}}$  that is independent of  $z_t$ , an appropriate  $\hat{\mathbf{y}}$  must be selected for the model input in practice. We experiment with two options for the independent condition. First,  $\hat{\mathbf{y}}$  can be drawn from a Gaussian distribution with a suitable standard deviation so that  $\hat{\mathbf{y}}$  matches the scale of the actual conditioning vector  $\mathbf{y}$ . Second, a random condition from the conditioning space, such as a random class label or random clip tokens, can be chosen as the independent  $\hat{\mathbf{y}}$ . We show in Section 7 that both methods perform similarly. However, there may be a slight preference for the random condition over Gaussian noise, as it stays closer to the conditioning distribution that the diffusion model was trained on.

## 5 TIME-STEP GUIDANCE

Inspired by ICG, we next offer an extension of classifier-free guidance that can be used with any model, including unconditional networks. We begin our analysis with class-conditional models and subsequently extend it to a more general setting.

In the class-conditional case, the embedding vector of the class is typically added to the embedding vector of the time step  $t$  to compute the input condition of the diffusion network. Hence, in practice, CFG essentially uses the outputs of the diffusion network for two different input embeddings and takes their difference as the update direction. Thus, we might directly utilize the time-step embedding of each diffusion model as a means to define a similar guidance signal. This leads to a novel method that we refer to as *time-step guidance* (TSG), which, like CFG, increases the quality of generations but, unlike CFG, is applicable even to unconditional models.

In this method, we compute the model outputs for the clean time-step embedding and a perturbed embedding and use their difference to guide the sampling. More specifically, at each time step  $t$ , we update the output via

$$\hat{D}_{\theta}(z_t, t) = D_{\theta}(z_t, \tilde{t}) + w_{\text{TSG}}(D_{\theta}(z_t, t) - D_{\theta}(z_t, \tilde{t})), \quad (8)$$

where  $\tilde{t}$  is the perturbed version of  $t$ . The intuition behind TSG is that at each time step  $t$ , altering the time-step embedding of the network leads to denoised outputs with either insufficient or excessive

noise removal (see Figure 10 in the appendix). Consequently, these outputs can be exploited to prevent the network from going toward undesirable predictions, thus increasing the accuracy of the score predictions at each time step. As we show below, TSG is related to stochastic Langevin dynamics in terms of the first-order approximation, and hence, is expected to improve generation quality.

**Connection to Langevin dynamics** Let  $\tilde{t} = t + \Delta t$ , where  $\Delta t$  is a small perturbation. Using a Taylor expansion, we get  $D_{\theta}(\mathbf{z}_t, \tilde{t}) = D_{\theta}(\mathbf{z}_t, t) + \frac{\partial D_{\theta}(\mathbf{z}_t, t)}{\partial t} \Delta t$ . Hence,  $\hat{D}_{\theta}(\mathbf{z}_t, \tilde{t}) = D_{\theta}(\mathbf{z}_t, t) + (1 - w_{\text{TSG}}) \frac{\partial D_{\theta}(\mathbf{z}_t, t)}{\partial t} \Delta t$ . Based on Equation (4), the score function is equal to

$$\nabla_{\mathbf{z}_t} \log \hat{p}_t(\mathbf{z}_t) = \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) + \frac{1 - w_{\text{TSG}}}{\sigma(t)^2} \frac{\partial D_{\theta}(\mathbf{z}_t, t)}{\partial t} \Delta t. \quad (9)$$

Now, if we follow the Euler sampling step for solving Equation (1), i.e. we define the update rule as  $\mathbf{z}_{t-1} = \mathbf{z}_t + \eta_t \nabla_{\mathbf{z}_t} \log \hat{p}_t(\mathbf{z}_t)$ , then the modified sampling step after time-step guidance will be equal to

$$\mathbf{z}_{t-1} = \mathbf{z}_t + \eta_t \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t) + \eta_t \frac{1 - w_{\text{TSG}}}{\sigma(t)^2} \frac{\partial D_{\theta}(\mathbf{z}_t, t)}{\partial t} \Delta t. \quad (10)$$

Assuming that  $\Delta t$  is a Gaussian random variable with zero mean, the update rule resembles a Langevin dynamics step, where the noise strength is determined based on the network behavior as represented by  $\frac{\partial D_{\theta}(\mathbf{z}_t, t)}{\partial t}$ . As Langevin dynamics is known to increase the quality of sampling from a given distribution by compensating for the errors happening at each sampling step, we argue that TSG also behaves similarly in terms of first-order approximation.

**Implementation details** In practice, we implement TSG by perturbing the time-step embedding with zero-mean Gaussian noise according to  $\tilde{t}_{\text{emb}} = t_{\text{emb}} + st^{\alpha} \mathbf{n}$  where  $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $st^{\alpha}$  determines the noise scale at each time step  $t$ . We choose  $s$  and  $\alpha$  such that the scale of the noise portion becomes comparable to the scale of the time-step embedding  $t_{\text{emb}}$ . Empirically, we also find that it is sometimes beneficial to apply the perturbed embeddings only to a portion of layers in the diffusion network, e.g., using  $\tilde{t}_{\text{emb}}$  for the first 10 layers and  $t_{\text{emb}}$  for the rest of layers. We provide ablations on these hyperparameters in Section 7.

## 6 EXPERIMENTS

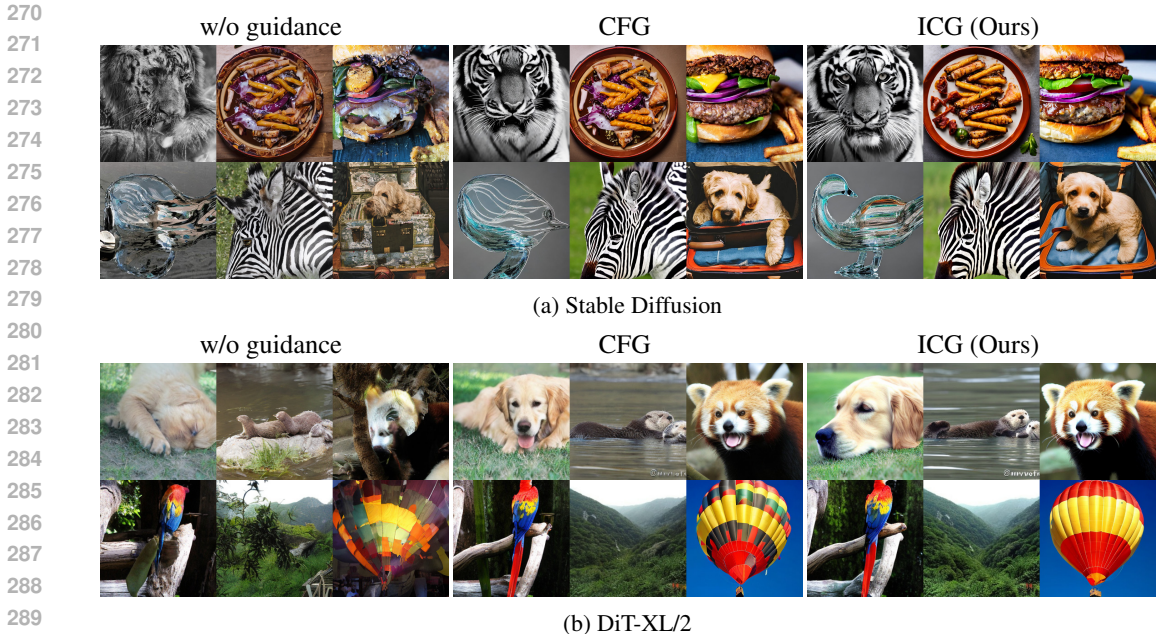
In this section, we rigorously evaluate ICG and demonstrate its ability to simulate the behavior of CFG across several conditional models. Additionally, we show that TSG improves the quality of both conditional and unconditional generations compared to the non-guided sampling baseline.

**Setup** All experiments are conducted via pre-trained checkpoints provided by official implementations. We use the recommended sampler that comes with each model, such as the EDM sampler for EDM networks (Karras et al., 2022), DPM++ (Lu et al., 2022b) for Stable Diffusion (Rombach et al., 2022), and DDPM (Ho et al., 2020) for DiT-XL/2 (Peebles & Xie, 2022).

**Evaluation** We use Fréchet Inception Distance (FID) (Heusel et al., 2017) as the main metric to measure both quality and diversity due to its alignment with human judgment. As FID is known to be sensitive to small implementation details, we ensure that models under comparison follow the same evaluation setup. For completeness, we also report precision (Kynkäänniemi et al., 2019) as a standalone quality metric and recall (Kynkäänniemi et al., 2019) as a diversity metric whenever possible.  $\text{FD}_{\text{DINOv2}}$  (Stein et al., 2024) is also reported for the EDM2 model (Karras et al., 2023).

### 6.1 COMPARISON BETWEEN ICG AND CFG

**Qualitative results** The qualitative comparisons between ICG and CFG are given in Figure 1 for Stable Diffusion (Rombach et al., 2022) and DiT-XL/2 (Peebles & Xie, 2022) models. Figure 2 also shows a comparison between the EDM2 model (Karras et al., 2023) guided with a separate unconditional module vs ICG. We observe that both ICG and CFG improve image quality, and the outputs of ICG and CFG are almost identical. This empirical evidence agrees with our theoretical justification provided in Section 4.



291 Figure 1: Comparison between CFG and ICG for (a) Stable Diffusion (Rombach et al., 2022) and (b)  
292 DiT-XL/2 (Peebles & Xie, 2022). Both CFG and ICG significantly improve the image quality of the  
293 baseline. Also note the similarity between the outputs of CFG and ICG, confirming our theoretical  
294 analysis in Section 4.



304 Figure 2: Comparison between the EDM2 model (Karras et al., 2023) guided with another unconditional  
305 module and ICG. We observe that using ICG leads to similar generations as CFG, and both  
306 methods significantly improve output quality compared to sampling without guidance.

307  
308  
309 **Quantitative results** We now show that ICG and CFG both result in similar performance metrics  
310 across several conditional models. As shown in Table 1, compared to CFG, ICG achieves better or  
311 similar performance across all metrics.<sup>1</sup> In Table 2, we also compare the effect of ICG on EDM  
312 (Karras et al., 2022) and EDM2 (Karras et al., 2023) models that were not trained with the CFG  
313 objective. The table shows that ICG performs similarly to guiding the generations with a separate  
314 unconditional module.

315 **Effect of removing the CFG objective from training** In this experiment, we demonstrate that the  
316 training component allocated to the CFG objective (i.e., label dropping) is unnecessary, and better  
317 results are obtained by training a purely conditional model and guiding the generations at inference  
318 with ICG. Using a DiT model for class-conditional ImageNet generation, Figure 3 shows that the  
319 purely conditional model consistently outperforms standard training with label dropping ( $p = 0.1$ )  
320 across all checkpoints. Consequently, training resources can be reallocated to the conditional part,  
321 leading to either faster convergence (by approximately 30%) or a better model (with around a 20%  
322 reduction in FID) with the same number of training iterations.

323 <sup>1</sup>For MDM (Tevet et al., 2023), recall is not available, and R-precision is reported similar to the paper.

Table 1: Quantitative comparison between CFG and ICG. ICG is able to achieve similar metrics to standard CFG by extracting the unconditional score from the conditional model itself.

Model	Architecture	Guidance	FID ↓	Precision ↑	Recall ↑
Stable Diffusion (Rombach et al., 2022)	UNet	CFG	20.13	<b>0.69</b>	<b>0.54</b>
		ICG (Ours)	<b>20.05</b>	<b>0.69</b>	0.53
DiT-XL/2 (Peebles & Xie, 2022)	Transformer	CFG	5.56	0.81	<b>0.66</b>
		ICG (Ours)	<b>5.50</b>	<b>0.83</b>	0.65
Pose-to-Image (Sadat et al., 2024)	UNet	CFG	14.61	0.93	0.02
		ICG (Ours)	<b>13.46</b>	<b>0.94</b>	<b>0.03</b>
MDM (Tevet et al., 2023)	Transformer	CFG	0.65	<b>0.73</b>	-
		ICG (Ours)	<b>0.47</b>	0.71	-

Table 2: Quantitative comparison between CFG and ICG for EDM networks. Although these models are not trained with the CFG objective, guiding their generations using a separate unconditional module results in similar outcomes to using ICG.

Model	Dataset	Guidance	FID ↓	FD <sub>DINOv2</sub> ↓
EDM2-XS (Karras et al., 2023)	Imagenet	CFG	3.36	79.94
		ICG (Ours)	<b>3.35</b>	<b>79.54</b>
EDM (Karras et al., 2022)	CIFAR-10	CFG	<b>1.87</b>	-
		ICG (Ours)	<b>1.87</b>	-

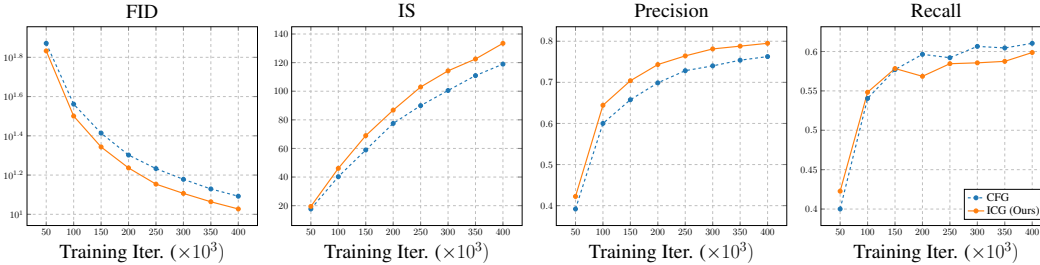


Figure 3: Comparison of CFG and ICG during training of a DiT model on ImageNet. Compared to standard CFG with label dropping, using ICG with a purely conditional model achieves better FID across all checkpoints. This indicates that the iterations spent on the CFG objective could be better allocated to training the conditional score, ultimately leading to a better model.

**Varying the Guidance Scale** Next, we demonstrate that by varying the guidance scale of ICG, we can increase the quality of outputs in a manner similar to standard CFG. As shown in Figure 4, increasing the guidance scale improves precision but reduces recall. The FID plots also form a U-shaped curve, consistent with what we expect from standard CFG.

**Results for ControlNet** We also show that ICG can be used for improving the quality of image-conditioned models as well. We use ControlNet (Zhang & Agrawala, 2023) as an example in this section since it is not trained with the CFG objective on the image condition input. That is, it only applies CFG to the text component of the condition. Our results are given in Figure 5. We see that without any text prompt, ICG significantly improves the quality of generations over the base sampling.

## 6.2 EFFECTIVENESS OF TIME-STEP GUIDANCE

Lastly, we show the effectiveness of time-step guidance in improving generation quality without relying on any information about the conditioning signal. The qualitative results are given in Figure 6. We can see that TSG increases the image quality of both conditional and unconditional sampling. Table 3 also presents the quantitative evaluation of TSG for both conditional and unconditional generation. Similar to CFG, using TSG significantly improves FID by trading diversity with quality. Finally, Figure 7 shows how TSG behaves as we increase the guidance scale. We observe that similar to CFG, TSG also has a U-shaped plot for the FID as the guidance scale increases.

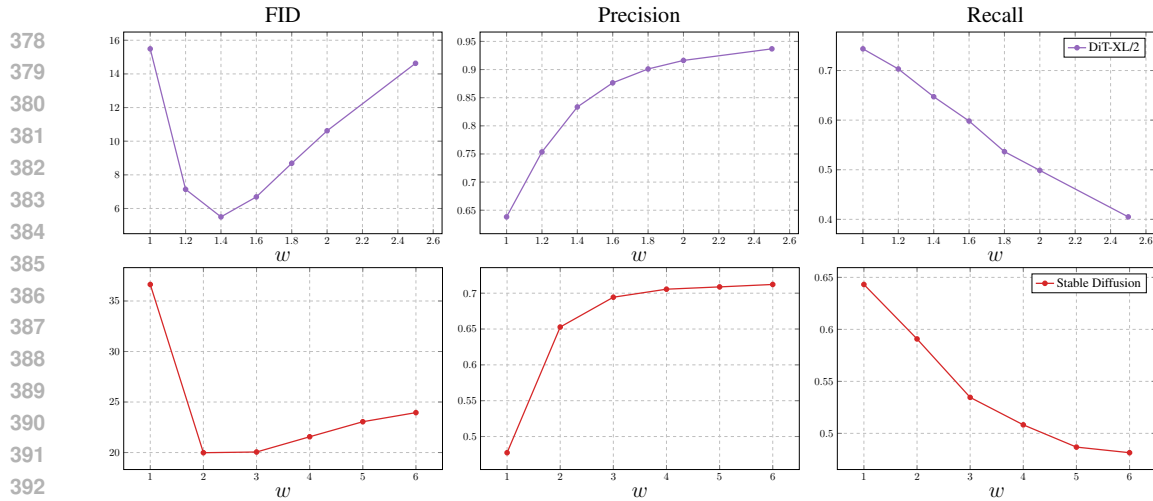


Figure 4: Behavior of ICG as the guidance scale increases. Similar to CFG, ICG trades diversity (lower recall) for quality (higher precision) at higher guidance scales.



(a) Pose-to-image generation



(b) Depth-to-image generation

Figure 5: Image-conditioned generation with ControlNet (without prompt). ICG significantly increases the quality of generations by applying guidance to the image condition.



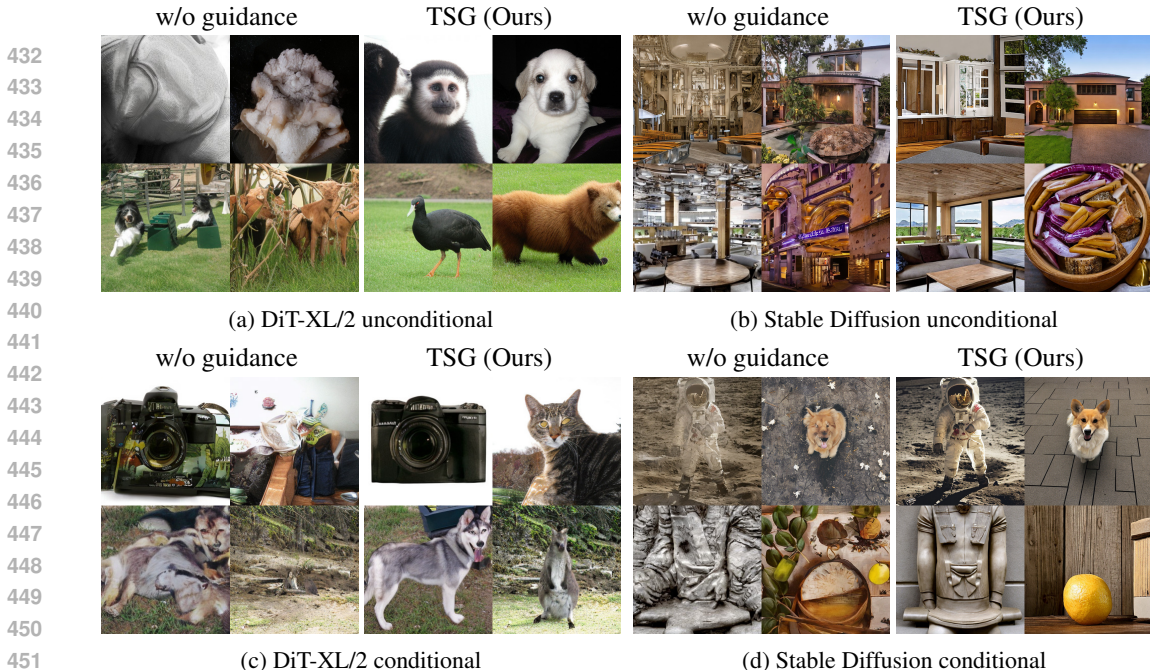


Figure 6: Effectiveness of TSG to improve the quality of both unconditional and conditional generation across two different models: DiT-XL/2 (Peebles & Xie, 2022) for class-conditional generation, and Stable Diffusion (Rombach et al., 2022) for text-conditional generation.

Table 3: Quantitative comparison between the baseline sampling of the diffusion models and sampling with TSG. TSG significantly boosts quality (lower FID) across various setups.

Model	Architecture	Type	Guidance	FID ↓	Precision ↑	Recall ↑
Stable Diffusion (Rombach et al., 2022)	UNet	Unconditional	✗	69.38	0.42	0.49
			TSG (Ours)	<b>56.65</b>	<b>0.54</b>	<b>0.54</b>
		Text-conditional	✗	36.63	0.48	<b>0.64</b>
			TSG (Ours)	<b>22.17</b>	<b>0.62</b>	0.59
DiT-XL/2 (Peebles & Xie, 2022)	Transformer	Unconditional	✗	48.67	0.48	<b>0.59</b>
			TSG (Ours)	<b>29.03</b>	<b>0.69</b>	0.55
		Class-conditional	✗	15.49	0.64	<b>0.74</b>
			TSG (Ours)	<b>6.39</b>	<b>0.82</b>	0.65

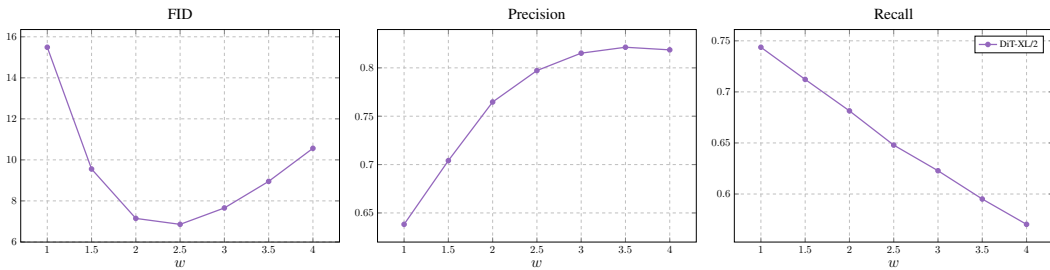


Figure 7: Behavior of TSG as the guidance scale increases for DiT-XL/2. Similar to CFG, TSG also significantly improves FID by trading diversity (recall) with quality (precision).

**Combining TSG and ICG** We also demonstrate that ICG and TSG can be complementary to each other when combined at the proper scale. The quantitative results of this experiment are presented in Table 4 with a visual example given in Figure 8. The table indicates that the combination of ICG and TSG outperforms each method in isolation in terms of FID, and all guided sampling algorithms significantly outperform the non-guided baseline.

Table 4: Compatibility of ICG and TSG

ICG	TSG	FID ↓	Precision ↑	Recall ↑
✗	✗	15.49	0.64	<b>0.74</b>
✓	✗	6.47	0.77	0.69
✗	✓	9.55	0.70	0.71
✓	✓	<b>5.76</b>	<b>0.82</b>	0.65

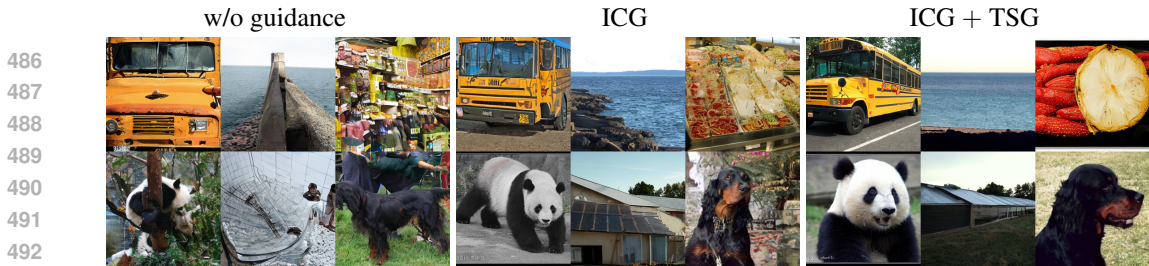


Figure 8: Visual example of combining ICG and TSG. The combination yields better visual generations compared to the baseline and using ICG alone.

Table 6: Ablation study examining various design elements in TSG.

(a) Influence of the noise scale $s$				(b) Influence of $\alpha$				(c) Maximum layer index			
$s$	FID ↓	Precision ↑	Recall ↑	$\alpha$	FID ↓	Precision ↑	Recall ↑	Index	FID ↓	Precision ↑	Recall ↑
1.0	10.23	0.69	0.35	0.75	7.22	0.82	0.62	5	7.84	0.75	<b>0.69</b>
2.0	<b>6.85</b>	<b>0.80</b>	<b>0.69</b>	1.0	<b>6.39</b>	<b>0.84</b>	0.65	10	<b>6.85</b>	0.79	0.65
2.5	7.94	<b>0.80</b>	<b>0.69</b>	1.25	6.47	0.78	<b>0.66</b>	15	7.65	<b>0.82</b>	0.65

## 7 ABLATION STUDIES

We next present the ablation studies on the effect of random conditioning in ICG and the hyperparameters in TSG. All ablations are conducted using the DiT-XL/2 model (Peebles & Xie, 2022).

**The choice of random condition in ICG** We first show that both Gaussian noise and random conditions can be used for estimating the unconditional part in ICG. The quantitative results are given in Table 5. The table shows that both methods are viable options for simulating classifier-free guidance without training.

Table 5: Ablation on the choice of independent condition in ICG.

ICG method	FID ↓	Precision ↑	Recall ↑
Gaussian noise	<b>5.50</b>	<b>0.83</b>	<b>0.65</b>
Random condition	5.55	<b>0.83</b>	<b>0.65</b>

**Impact of hyperparameters in time-step guidance** This ablation study explores the effect of hyperparameters in TSG. The results are presented in Table 6. We observe that as we introduce more perturbation into the time-step embedding of the model, in the form of higher noise scale  $s$  (Table 6a), lower power  $\alpha$  (Table 6b), or higher layer index (Table 6c), precision improves while recall decreases. This suggests that the amount of perturbation should be balanced for a good trade-off between diversity and quality. We also empirically observed that adding too much noise to the time-step embedding hurts image quality.

## 8 DISCUSSION AND CONCLUSION

In this paper, we revisited the core aspects of classifier-free guidance and showed that by replacing the conditional vector in a trained conditional diffusion model with an independent condition, we can efficiently estimate the score of the unconditional distribution. We then introduced independent condition guidance (ICG), a novel method that simulates the same behavior as CFG without the need to learn an unconditional model during training. Inspired by this, we also proposed time-step guidance (TSG) and demonstrated that the time-step information learned by the diffusion model can be leveraged to enhance the quality of generations, even for unconditional models. Our experiments indicate that ICG performs similarly to standard CFG and alleviates the need to consider the CFG objective during training. Thus, ICG streamlines the training of conditional models and improves training efficiency. Additionally, we verified that TSG also improves generation quality in a manner similar to CFG, without relying on any conditional information. As with CFG, challenges remain in accelerating the proposed methods to narrow the gap between the cost of guided and unguided sampling (i.e., eliminating the need to query the diffusion model twice per sampling step); we view this topic as a promising avenue for further research.

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

## ETHICS STATEMENT

As generative modeling advances, the creation and spread of fabricated or inaccurate data become easier. Thus, while improvements in AI-generated content can boost productivity and creativity, it is crucial to consider the associated risks and ethical implications. For a more detailed discussion on the ethics and creativity in computer vision, we refer readers to Rostamzadeh et al. (2021).

## REPRODUCIBILITY STATEMENT

This work is based on the official implementations of the pretrained models referenced in the main text. The exact algorithms for ICG and TSG are provided in Algorithms 1 and 2, with corresponding pseudocode shown in Figures 11 and 12. Additional implementation details, including the specific hyperparameters used to generate the results in this paper, are discussed in Appendix E.

## REFERENCES

- Donghoon Ahn, Hyoungwon Cho, Jaewon Min, Wooseok Jang, Jungwoo Kim, Seonhwa Kim, Hyun Hee Park, Kyong Hwan Jin, and Seungryong Kim. Self-rectifying diffusion sampling with perturbed-attention guidance. *CoRR*, abs/2403.17377, 2024. doi: 10.48550/ARXIV.2403.17377. URL <https://doi.org/10.48550/arXiv.2403.17377>.
- Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu Liu. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *CoRR*, abs/2211.01324, 2022. doi: 10.48550/arXiv.2211.01324. URL <https://doi.org/10.48550/arXiv.2211.01324>.
- Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 843–852, 2023.
- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. Stable video diffusion: Scaling latent video diffusion models to large datasets. *CoRR*, abs/2311.15127, 2023a. doi: 10.48550/ARXIV.2311.15127. URL <https://doi.org/10.48550/arXiv.2311.15127>.
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. 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, 2023b.
- Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=Blxsqj09Fm>.
- Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=NsMLjcFaO8O>.
- Hyungjin Chung, Jeongsol Kim, Michael T Mccann, Marc L Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. *arXiv preprint arXiv:2209.14687*, 2022.
- Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat gans on image synthesis. In Marc Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 8780–8794, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/49ad23d1ec9fa4bd8d77d02681df5cfa-Abstract.html>.

- 594 Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang,  
595 and José Lezama. Photorealistic video generation with diffusion models. *arXiv preprint*  
596 *arXiv:2312.06662*, 2023.
- 597 Yutong He, Naoki Murata, Chieh-Hsin Lai, Yuhta Takida, Toshimitsu Uesaka, Dongjun Kim, Wei-  
598 Hsiang Liao, Yuki Mitsufuji, J Zico Kolter, Ruslan Salakhutdinov, et al. Manifold preserving  
599 guided diffusion. *arXiv preprint arXiv:2311.16424*, 2023.
- 600 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
601 Gans trained by a two time-scale update rule converge to a local nash equilibrium. In Isabelle  
602 Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vish-  
603 wanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30:*  
604 *Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long*  
605 *Beach, CA, USA*, pp. 6626–6637, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/](https://proceedings.neurips.cc/paper/2017/hash/8a1d694707eb0fefe65871369074926d-Abstract.html)  
606 [8a1d694707eb0fefe65871369074926d-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/8a1d694707eb0fefe65871369074926d-Abstract.html).
- 607 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *CoRR*, abs/2207.12598, 2022.  
608 doi: 10.48550/arXiv.2207.12598. URL <https://doi.org/10.48550/arXiv.2207.12598>.
- 609 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo  
610 Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.),  
611 *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information*  
612 *Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL [https:](https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html)  
613 [/proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html).
- 614 Susung Hong, Gyuseong Lee, Wooseok Jang, and Seungryong Kim. Improving sample quality of  
615 diffusion models using self-attention guidance. *CoRR*, abs/2210.00939, 2022. doi: 10.48550/arXiv.  
616 2210.00939. URL <https://doi.org/10.48550/arXiv.2210.00939>.
- 617 Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for  
618 high resolution images. *CoRR*, abs/2301.11093, 2023. doi: 10.48550/arXiv.2301.11093. URL  
619 <https://doi.org/10.48550/arXiv.2301.11093>.
- 620 Qingqing Huang, Daniel S. Park, Tao Wang, Timo I. Denk, Andy Ly, Nanxin Chen, Zhengdong  
621 Zhang, Zhishuai Zhang, Jiahui Yu, Christian Havnø Frank, Jesse H. Engel, Quoc V. Le, William  
622 Chan, and Wei Han. Noise2music: Text-conditioned music generation with diffusion models.  
623 *CoRR*, abs/2302.03917, 2023. doi: 10.48550/arXiv.2302.03917. URL [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.2302.03917)  
624 [arXiv.2302.03917](https://doi.org/10.48550/arXiv.2302.03917).
- 625 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-  
626 based generative models. 2022. URL <https://openreview.net/forum?id=k7FuTOWMOc7>.
- 627 Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing  
628 and improving the training dynamics of diffusion models, 2023.
- 629 Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile  
630 diffusion model for audio synthesis. In *9th International Conference on Learning Representations,*  
631 *ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL [https://openreview.](https://openreview.net/forum?id=a-xFK8Ymz5J)  
632 [net/forum?id=a-xFK8Ymz5J](https://openreview.net/forum?id=a-xFK8Ymz5J).
- 633 Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Im-  
634 proved precision and recall metric for assessing generative models. In Hanna M. Wallach,  
635 Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Gar-  
636 nett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on*  
637 *Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancou-*  
638 *ver, BC, Canada*, pp. 3929–3938, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/0234c510bc6d908b28c70ff313743079-Abstract.html)  
639 [0234c510bc6d908b28c70ff313743079-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/0234c510bc6d908b28c70ff313743079-Abstract.html).
- 640 Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
641 Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In David J.  
642 Fleet, Tomás Pajdla, Bernt Schiele, and Tinne Tuytelaars (eds.), *Computer Vision - ECCV 2014*  
643 *- 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V,*  
644 *volume 8693 of Lecture Notes in Computer Science*, pp. 740–755. Springer, 2014. doi: 10.1007/  
645 978-3-319-10602-1\_48. URL [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-319-10602-1_48)  
646 [978-3-319-10602-1\\_48](https://doi.org/10.1007/978-3-319-10602-1_48).

- 648 Guan-Hong Liu, Arash Vahdat, De-An Huang, Evangelos A. Theodorou, Weili Nie, and Anima  
649 Anandkumar.  $I^2_{sb}$ : Image-to-image schrödinger bridge. *CoRR*, abs/2302.05872, 2023. doi:  
650 10.48550/arXiv.2302.05872. URL <https://doi.org/10.48550/arXiv.2302.05872>.  
651
- 652 Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models  
653 on manifolds. In *The Tenth International Conference on Learning Representations, ICLR 2022,*  
654 *Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=PIKWVd2yBkY)  
655 [PIKWVd2yBkY](https://openreview.net/forum?id=PIKWVd2yBkY).
- 656 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-  
657 solver: A fast ODE solver for diffusion probabilistic model sampling in around  
658 10 steps. In *NeurIPS, 2022a*. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/260a14acce2a89dad36adc8eefe7c59e-Abstract-Conference.html)  
659 [260a14acce2a89dad36adc8eefe7c59e-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/260a14acce2a89dad36adc8eefe7c59e-Abstract-Conference.html).  
660
- 661 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast  
662 solver for guided sampling of diffusion probabilistic models. *CoRR*, abs/2211.01095, 2022b. doi:  
663 10.48550/arXiv.2211.01095. URL <https://doi.org/10.48550/arXiv.2211.01095>.
- 664 Kevin P. Murphy. *Probabilistic Machine Learning: Advanced Topics*. MIT Press, 2023. URL  
665 <http://probml.github.io/book2>.  
666
- 667 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
668 In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on*  
669 *Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of*  
670 *Machine Learning Research*, pp. 8162–8171. PMLR, 2021. URL [http://proceedings.mlr.press/](http://proceedings.mlr.press/v139/nichol21a.html)  
671 [v139/nichol21a.html](http://proceedings.mlr.press/v139/nichol21a.html).
- 672 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob  
673 McGrew, Ilya Sutskever, and Mark Chen. GLIDE: towards photorealistic image generation and  
674 editing with text-guided diffusion models. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song,  
675 Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine*  
676 *Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of*  
677 *Machine Learning Research*, pp. 16784–16804. PMLR, 2022. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v162/nichol22a.html)  
678 [v162/nichol22a.html](https://proceedings.mlr.press/v162/nichol22a.html).
- 679 William Peebles and Saining Xie. Scalable diffusion models with transformers. *CoRR*,  
680 abs/2212.09748, 2022. doi: 10.48550/arXiv.2212.09748. URL [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2212.09748)  
681 [2212.09748](https://doi.org/10.48550/arXiv.2212.09748).  
682
- 683 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
684 Penna, and Robin Rombach. SDXL: improving latent diffusion models for high-resolution  
685 image synthesis. *CoRR*, abs/2307.01952, 2023. doi: 10.48550/ARXIV.2307.01952. URL [https://](https://doi.org/10.48550/arXiv.2307.01952)  
686 [doi.org/10.48550/arXiv.2307.01952](https://doi.org/10.48550/arXiv.2307.01952).
- 687 Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d  
688 diffusion. In *The Eleventh International Conference on Learning Representations, ICLR 2023,*  
689 *Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.net/pdf?id=](https://openreview.net/pdf?id=FjNys5c7VyY)  
690 [FjNys5c7VyY](https://openreview.net/pdf?id=FjNys5c7VyY).  
691
- 692 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-  
693 conditional image generation with CLIP latents. *CoRR*, abs/2204.06125, 2022. doi: 10.48550/  
694 arXiv.2204.06125. URL <https://doi.org/10.48550/arXiv.2204.06125>.
- 695 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
696 resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer*  
697 *Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp.  
698 10674–10685. IEEE, 2022. doi: 10.1109/CVPR52688.2022.01042. URL [https://doi.org/10.1109/](https://doi.org/10.1109/CVPR52688.2022.01042)  
699 [CVPR52688.2022.01042](https://doi.org/10.1109/CVPR52688.2022.01042).
- 700  
701 Negar Rostamzadeh, Emily Denton, and Linda Petrini. Ethics and creativity in computer vision.  
*CoRR*, abs/2112.03111, 2021. URL <https://arxiv.org/abs/2112.03111>.

- 702 Seyedmorteza Sadat, Jakob Buhmann, Derek Bradley, Otmar Hilliges, and Romann M. Weber. CADs:  
703 Unleashing the diversity of diffusion models through condition-annealed sampling. In *The Twelfth*  
704 *International Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=zMoNrajK2X)  
705 [id=zMoNrajK2X](https://openreview.net/forum?id=zMoNrajK2X).
- 706 Chitwan Saharia, William Chan, Huiwen Chang, Chris A. Lee, Jonathan Ho, Tim Salimans, David J.  
707 Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In Munkhtsetseg  
708 Nandigjav, Niloy J. Mitra, and Aaron Hertzmann (eds.), *SIGGRAPH '22: Special Interest Group*  
709 *on Computer Graphics and Interactive Techniques Conference, Vancouver, BC, Canada, August 7 -*  
710 *11, 2022*, pp. 15:1–15:10. ACM, 2022a. doi: 10.1145/3528233.3530757. URL [https://doi.org/10.](https://doi.org/10.1145/3528233.3530757)  
711 [1145/3528233.3530757](https://doi.org/10.1145/3528233.3530757).
- 712 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kam-  
713 yar Seyed Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan  
714 Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with  
715 deep language understanding. 2022b. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/ec795aeadae0b7d230fa35cbaf04c041-Abstract-Conference.html)  
716 [ec795aeadae0b7d230fa35cbaf04c041-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/ec795aeadae0b7d230fa35cbaf04c041-Abstract-Conference.html).
- 717 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In  
718 *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April*  
719 *25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=TIIdIXIpzhoI>.
- 720 Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
721 learning using nonequilibrium thermodynamics. *37:2256–2265*, 2015. URL [http://proceedings.](http://proceedings.mlr.press/v37/sohl-dickstein15.html)  
722 [mlr.press/v37/sohl-dickstein15.html](http://proceedings.mlr.press/v37/sohl-dickstein15.html).
- 723 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *9th*  
724 *International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May*  
725 *3-7, 2021*. OpenReview.net, 2021a. URL <https://openreview.net/forum?id=St1giarCHLP>.
- 726 Jiaming Song, Qinsheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin  
727 Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation.  
728 In *International Conference on Machine Learning*, pp. 32483–32498. PMLR, 2023.
- 729 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.  
730 In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B.  
731 Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual*  
732 *Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019,*  
733 *Vancouver, BC, Canada*, pp. 11895–11907, 2019. URL [https://proceedings.neurips.cc/paper/2019/](https://proceedings.neurips.cc/paper/2019/hash/3001ef257407d5a371a96dcd947c7d93-Abstract.html)  
734 [hash/3001ef257407d5a371a96dcd947c7d93-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/3001ef257407d5a371a96dcd947c7d93-Abstract.html).
- 735 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and  
736 Ben Poole. Score-based generative modeling through stochastic differential equations. In *9th*  
737 *International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May*  
738 *3-7, 2021*. OpenReview.net, 2021b. URL <https://openreview.net/forum?id=PxTIG12RRHS>.
- 739 George Stein, Jesse Cresswell, Rasa Hosseinzadeh, Yi Sui, Brendan Ross, Valentin Vilecroze,  
740 Zhaoyan Liu, Anthony L Caterini, Eric Taylor, and Gabriel Loaiza-Ganem. Exposing flaws of  
741 generative model evaluation metrics and their unfair treatment of diffusion models. *Advances in*  
742 *Neural Information Processing Systems*, 36, 2024.
- 743 Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit Haim Bermano.  
744 Human motion diffusion model. 2023. URL <https://openreview.net/pdf?id=SJ1kSyO2jwu>.
- 745 Jonathan Tseng, Rodrigo Castellon, and C. Karen Liu. EDGE: editable dance generation from music.  
746 In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver,*  
747 *BC, Canada, June 17-24, 2023*, pp. 448–458. IEEE, 2023. doi: 10.1109/CVPR52729.2023.00051.  
748 URL <https://doi.org/10.1109/CVPR52729.2023.00051>.
- 749 Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan,  
750 Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin  
751 Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich  
752 text-to-image generation. *Trans. Mach. Learn. Res.*, 2022, 2022. URL [https://openreview.net/](https://openreview.net/forum?id=AFDcYJKhND)  
753 [forum?id=AFDcYJKhND](https://openreview.net/forum?id=AFDcYJKhND).

756 Jiwen Yu, Yinhuai Wang, Chen Zhao, Bernard Ghanem, and Jian Zhang. Freedom: Training-  
757 free energy-guided conditional diffusion model. In *Proceedings of the IEEE/CVF International*  
758 *Conference on Computer Vision*, pp. 23174–23184, 2023.

759  
760 Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models.  
761 *CoRR*, abs/2302.05543, 2023. doi: 10.48550/ARXIV.2302.05543. URL [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.2302.05543)  
762 [arXiv.2302.05543](https://doi.org/10.48550/arXiv.2302.05543).

763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## A SUFFICIENCY OF THE CONDITIONAL SCORE

In this section, we provide another perspective on why training the conditional score is sufficient for computing the unconditional score. For ease of exposition, assume that the model’s objective is to directly learn the conditional score,  $s_\theta(\mathbf{x}, \mathbf{y}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{y})$ , from paired data jointly drawn from  $p(\mathbf{x}, \mathbf{y})$ . This can be achieved by directly using denoising score matching (Song et al., 2021b) or by training a denoiser  $D_\theta(\mathbf{x}, \mathbf{y})$  and converting its outputs to the corresponding score function via Equation (4). The question then boils down to whether the unconditional score,  $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ , can be recovered from the conditional score. By direct calculation, we have that

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}) = \frac{\nabla_{\mathbf{x}} p(\mathbf{x})}{p(\mathbf{x})} \quad (11)$$

$$= \frac{\nabla_{\mathbf{x}} \int p(\mathbf{y}) p(\mathbf{x}|\mathbf{y}) d\mathbf{y}}{p(\mathbf{x})} \quad (12)$$

$$= \int \frac{p(\mathbf{y})}{p(\mathbf{x})} \nabla_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) d\mathbf{y} \quad (13)$$

$$= \int \frac{p(\mathbf{y}) p(\mathbf{x}|\mathbf{y})}{p(\mathbf{x}) p(\mathbf{x}|\mathbf{y})} \nabla_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) d\mathbf{y} \quad (14)$$

$$= \int p(\mathbf{y}|\mathbf{x}) \nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{y}) d\mathbf{y}. \quad (15)$$

The last term is the conditional expectation of the conditional score, which is implicitly learned during the training of  $s_\theta(\mathbf{x}, \mathbf{y})$ . This shows that under sufficient data and optimization, the unconditional score at each time step is theoretically available from the conditional score through (conditional) marginalization.

In Section 4, we rely on the simpler relation given in Equation (7), which shows that when the condition is independent of the main argument of the score, the resulting “independent-conditional” score is equivalent to the unconditional score.

## B COMPATIBILITY OF ICG WITH CADS

We show that ICG is compatible with CADS (Sadat et al., 2024), and CADS can be used on top of ICG to increase the diversity of generations. An example of applying CADS to ICG is shown in Figure 9, and the quantitative results are given in Table 7 for the DiT-XL/2 model. As ICG behaves similarly to the standard CFG, applying CADS increases diversity with minimal drop in quality.

Table 7: Effectiveness of CADS on ICG.

Guidance	FID ↓	Precision ↑	Recall ↑
ICG	20.56	<b>0.89</b>	0.32
+CADS	<b>8.83</b>	0.78	<b>0.61</b>

## C INTUITION BEHIND TSG

This section provides more intuition on time-step guidance. We demonstrate that if we perturb the time step with a positive or negative constant (using  $t + \delta$  or  $t - \delta$ ) to guide the sampling, it results in insufficient or excessive noise removal in final generations. As shown in Figure 10, using lower time steps causes the model to perform excessive noise removal (soft outputs), while using higher time steps forces the model to perform insufficient noise removal (noisy images). TSG uses both directions to prevent the outputs from moving toward these undesirable paths, thereby increasing the quality of the generations.

## D INCREASING THE NUMBER OF SAMPLING STEPS

Since TSG increases the number of sampling steps (similar to CFG), a natural question is whether the same behavior can be achieved by simply increasing the number of sampling steps in the unguided sampling baseline. Our results in Table 8 indicate that this approach performs significantly worse than TSG, suggesting that, like CFG, TSG alters the sampling trajectories toward higher-quality



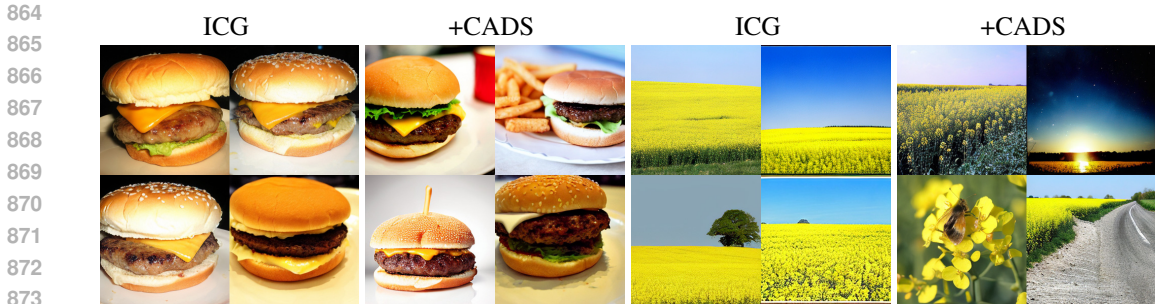


Figure 9: Similar to CFG, ICG is compatible with CADS, and CADS can be used to increase the diversity of ICG at higher guidance scales. Samples are generated from the DiT-XL/2 model.

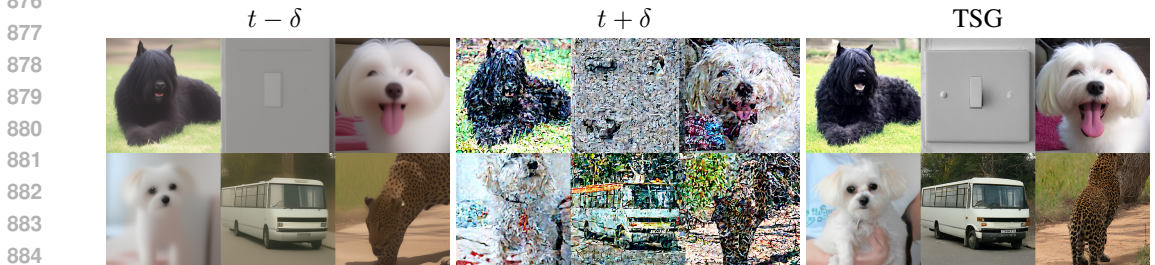


Figure 10: Intuition behind TSG: Using lower time steps for guidance causes excessive noise removal (soft outputs), while higher time steps cause insufficient noise removal (noisy images). TSG employs both directions to improve output quality.

Table 8: Comparison between unguided sampling and TSG using different number of sampling steps. TSG is able to achieve significantly better FID compared to both unguided baselines.

Guidance	# Steps	FID	Precision	Recall
Unguided	100	15.49	0.6382	<b>0.7437</b>
Unguided	200	12.94	0.6683	0.7387
TSG	100	<b>6.39</b>	<b>0.8198</b>	0.6489

generations. Note that TSG achieves a significantly better FID compared to both unguided sampling baselines. This aligns with findings from the CFG literature, where guided sampling outperforms unguided baselines even with many sampling steps (e.g., 1000 steps) (Dhariwal & Nichol, 2021; Ho & Salimans, 2022).

## E IMPLEMENTATION DETAILS

The sampling details for ICG and TSG are provided in Algorithms 1 and 2. Both algorithms are straightforward to implement and require minimal code changes compared to the base sampling. The pseudocode for implementing ICG and TSG is also included in Figures 11 and 12. Additionally, the hyperparameters used in our experiments are listed in Tables 9 and 10. The CADS experiment was conducted with a linear schedule using  $\tau_1 = 0.5$ ,  $\tau_2 = 0.9$ , and  $s_{CADS} = 0.15$ . Lastly, please note that we did not perform an exhaustive grid search on the parameters of TSG, and better configurations are likely to exist for each model.

For the TSG noise schedule, we experimented with a constant and a power schedule, as shown in Figure 12, and found that both work similarly. We recommend using the power schedule as it offers more flexibility over the scale of the noise at each  $t$ . The constant schedule is technically a special case of the power schedule, where the exponent is zero. We also found it useful to apply TSG only at intervals during the sampling, i.e., for  $t \in [T_{min}, T_{max}]$ , where  $T_{min}$  and  $T_{max}$  are hyperparameters. Also, when limiting TSG to only a portion of layers in the diffusion model, we used the first  $N$  layers of transformer-based architectures and the first  $N$  layers of the encoder and decoder in UNet-based

918 architectures. We chose to scale the amount of noise  $s$  based on the standard deviation of the time-step  
919 embedding (see Figures 11 and 12) for more fine-grained control over the scale.

920  
921 We primarily use the ADM evaluation script (Dhariwal & Nichol, 2021) for computing FID, precision,  
922 and recall to ensure a fair comparison across experiments. For class-conditional models, the FID is  
923 computed between 10,000 (for DiT-XL/2) or 50,000 (For EDM and EDM2) generated images and  
924 the full training dataset. For text-to-image models, we use the evaluation subset of MS COCO 2017  
925 (Lin et al., 2014) as the ground truth for captions and images.

## 926 F MORE VISUAL RESULTS

927  
928 This section presents additional visual results for our guidance methods. More results on ICG are  
929 provided in Figure 13, while additional results for TSG are shown in Figures 14 and 15. Consistent  
930 with the main results of the paper, ICG produces similar outcomes to CFG, and TSG consistently  
931 enhances the quality compared to the baseline. Figure 16 provides examples of the effectiveness  
932 of TSG based on Stable Diffusion XL (SDXL) (Podell et al., 2023). Finally, Figure 17 shows  
933 a qualitative comparison between unguided sampling and sampling with TSG for several latent  
934 diffusion models from Rombach et al. (2022).  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

---

**Algorithm 1** Sampling with ICG

---

**Require:**  $w_{\text{ICG}}$ : ICG strength  
**Require:**  $\mathbf{y}$ : Input condition  
 1: Initial value:  $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
 2: **for**  $t = T, \dots, 1$  **do**  
   3:   ○ Pick a random  $\hat{\mathbf{y}}$  independent of  $\mathbf{z}_t$  (Gaussian noise or from the conditioning space).  
   4:   ○ Compute the ICG guided output at  $t$ :  
        $\hat{D}_{\text{ICG}}(\mathbf{z}_t, t, \mathbf{y}) = D(\mathbf{z}_t, t, \hat{\mathbf{y}}) + w_{\text{ICG}}(D(\mathbf{z}_t, t, \mathbf{y}) - D(\mathbf{z}_t, t, \hat{\mathbf{y}}))$ .  
   5:   ○ Perform one sampling step (e.g. one step of DDPM):  
        $\mathbf{z}_{t-1} = \text{diffusion\_reverse}(\hat{D}_{\text{ICG}}, \mathbf{z}_t, t)$ .  
 6: **end for**  
 7: **return**  $\mathbf{z}_0$

---



---

**Algorithm 2** Sampling with TSG

---

**Require:**  $w_{\text{TSG}}$ : TSG strength  
**Require:**  $(s, \alpha)$ : TSG hyperparameters  
**Require:**  $\mathbf{y}$ : Input condition (optional)  
 1: Initial value:  $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
 2: **for**  $t = T, \dots, 1$  **do**  
   3:   ○ Perturb the time-step embedding  $t_{\text{emb}}$  to get  $\hat{t}_{\text{emb}}$ :  
        $\hat{t}_{\text{emb}} = t_{\text{emb}} + st^\alpha \mathbf{n}$ , where  $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .  
   4:   ○ Compute the TSG guided output at  $t$ :  
        $\hat{D}_{\text{TSG}}(\mathbf{z}_t, t, \mathbf{y}) = D(\mathbf{z}_t, \hat{t}_{\text{emb}}, \mathbf{y}) + w_{\text{TSG}}(D(\mathbf{z}_t, t_{\text{emb}}, \mathbf{y}) - D(\mathbf{z}_t, \hat{t}_{\text{emb}}, \mathbf{y}))$ .  
   5:   ○ Perform one sampling step (e.g. one step of DDPM):  
        $\mathbf{z}_{t-1} = \text{diffusion\_reverse}(\hat{D}_{\text{TSG}}, \mathbf{z}_t, t)$ .  
 6: **end for**  
 7: **return**  $\mathbf{z}_0$

---

Table 9: Hyperparameters used for the ICG experiments.

Model	ICG mode	ICG scale	CFG scale
DiT-XL/2	Random class	1.4	1.5
Stable Diffusion	Random text	3.0	4.0
Pose-to-Image	Gaussian noise	3.0	4.0
MDM	Gaussian noise	2.5	2.5
EDM	Random class	1.05	1.1
EDM2	Random class	1.25	1.25

Table 10: Hyperparameters used for the TSG experiments.

Model	Mode	TSG function	TSG scale	TSG parameters
DiT-XL/2	Unconditional	constant_schedule	5.0	T_MIN = 200, T_MAX = 800, s = 1.0
DiT-XL/2	Conditional	power_schedule	2.5	T_MIN = 0, T_MAX = 1000, $\alpha = 1$ , s = 2
Stable Diffusion	Unconditional	constant_schedule	3.0	T_MIN = 100, T_MAX = 900, s = 1.25
Stable Diffusion	Conditional	power_schedule	4.0	T_MIN = 400, T_MAX = 1000, s = 3.0, $\alpha = 0.25$

```

1026
1027
1028 1 def get_random_class():
1029 2     """Random class labels."""
1030 3     y_random = torch.randint(0, NUM_CLASSES, (BATCH_SIZE, ))
1031 4     return y_random
1032
1033 5
1034 6 def get_random_text():
1035 7     """Random text tokens."""
1036 8     random_idx = torch.randint(0, NUM_TOKENS, (BATCH_SIZE, MAX_LENGTH))
1037 9     random_tokens = text_encoder(random_idx, attention_mask=None)[0]
1038 10    return random_tokens
1039
1040 11
1041 12 def get_gaussian_noise_embedding(embeddings):
1042 13     """Random embedding based on Gaussian noise."""
1043 14     noise_embedding = torch.randn_like(embeddings) * embeddings.std()
1044 15     return noise_embedding
1045
1046 16
1047 17 def get_gaussian_noise_image(image_cond):
1048 18     """Random condition for image-conditional models."""
1049 19     noise_embedding = torch.randn_like(image_cond) * SCALE
1050 20     return noise_embedding
1051
1052
1053

```

Figure 11: Implementation details for ICG. The figure presents pseudocode for implementing the random class, random text, and Gaussian noise embedding for the unconditional component in ICG.

```

1054 1 def get_constant_schedule(t_emb, t, std_scaling=True):
1055 2     """Applies TSG for a portion of sampling (t in [T_MIN, T_MAX])."""
1056 3     if t < T_MIN or t > T_MAX:
1057 4         return t_emb
1058 5
1059 6     noise_scale = S
1060 7     if std_scaling:
1061 8         noise_scale = S * t_emb.std()
1062 9     that_emb = t_emb + torch.randn_like(t_emb) * noise_scale
1063 10    return that_emb
1064
1065
1066 11
1067 12 def get_power_schedule(t_emb, t, std_scaling=True):
1068 13     """Applies TSG according to the power schedule."""
1069 14     if t < T_MIN or t > T_MAX:
1070 15         return t_emb
1071 16     noise_scale = S * t ** (ALPHA)
1072 17     if std_scaling:
1073 18         noise_scale = noise_scale * t_emb.std()
1074 19     that_emb = t_emb + torch.randn_like(t_emb) * noise_scale
1075 20     return that_emb
1076
1077

```

Figure 12: Implementation details for TSG. We provide two scheduling techniques for perturbing the time-step embedding. We empirically found that both methods perform similarly.

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133



(a) Stable Diffusion



(b) DiT-XL/2

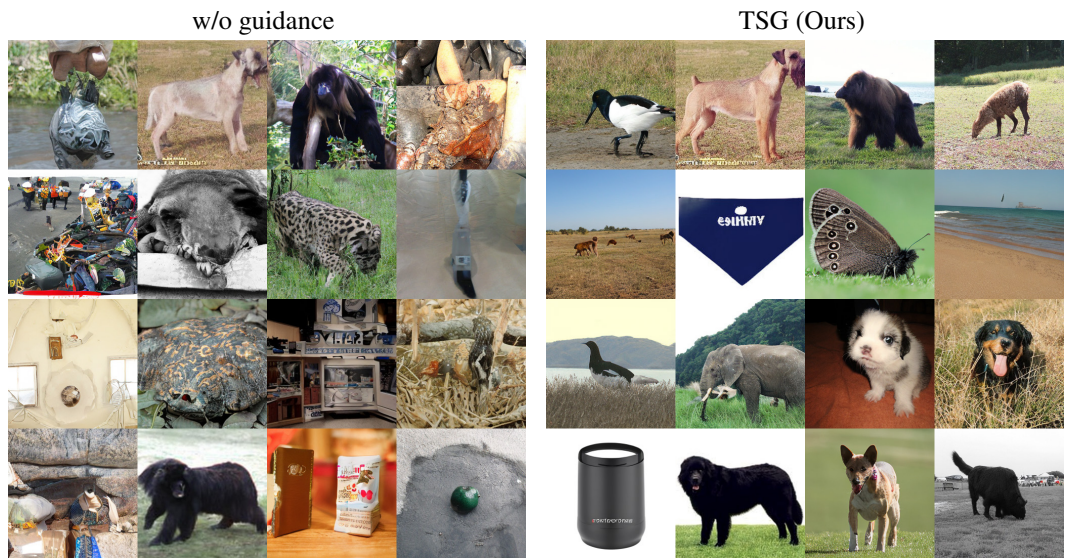
Figure 13: More comparisons between ICG and CFG for (a) text-to-image generation with Stable Diffusion (Rombach et al., 2022) and (b) class-conditional generation with DiT-XL/2 (Peebles & Xie, 2022).

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

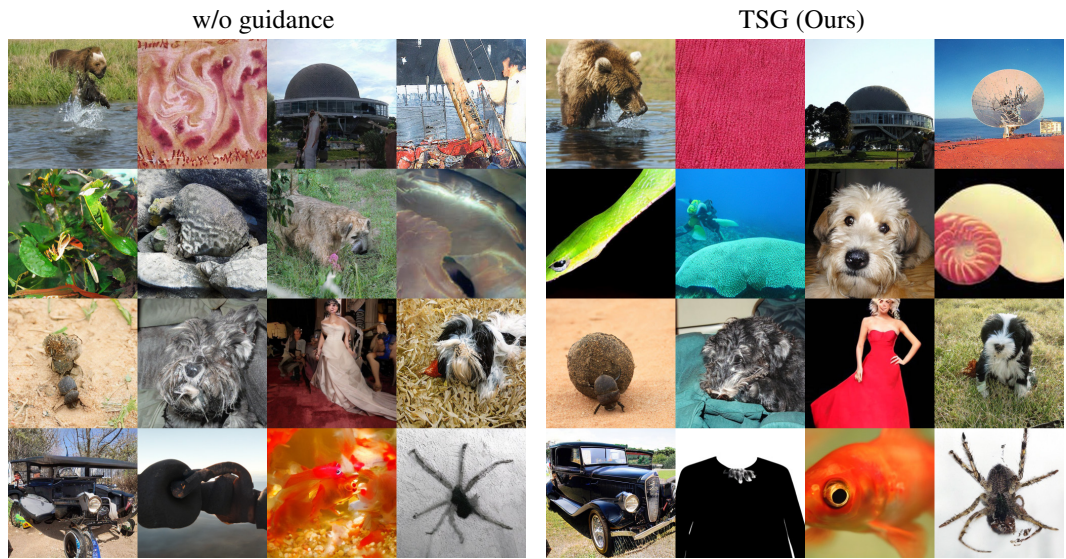


Figure 14: More comparisons on the effectiveness of TSG for improving the quality of both unconditional and conditional generation for Stable Diffusion (Rombach et al., 2022).

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241



(a) DiT-XL/2 unconditional



(b) DiT-XL/2 conditional

Figure 15: More comparisons on the effectiveness of TSG for improving the quality of both unconditional and conditional generation for DiT-XL/2 (Peebles & Xie, 2022).



Figure 16: Showcasing the effectiveness of TSG in improving the quality of generations compared to sampling without guidance based on SDXL (Podell et al., 2023).



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

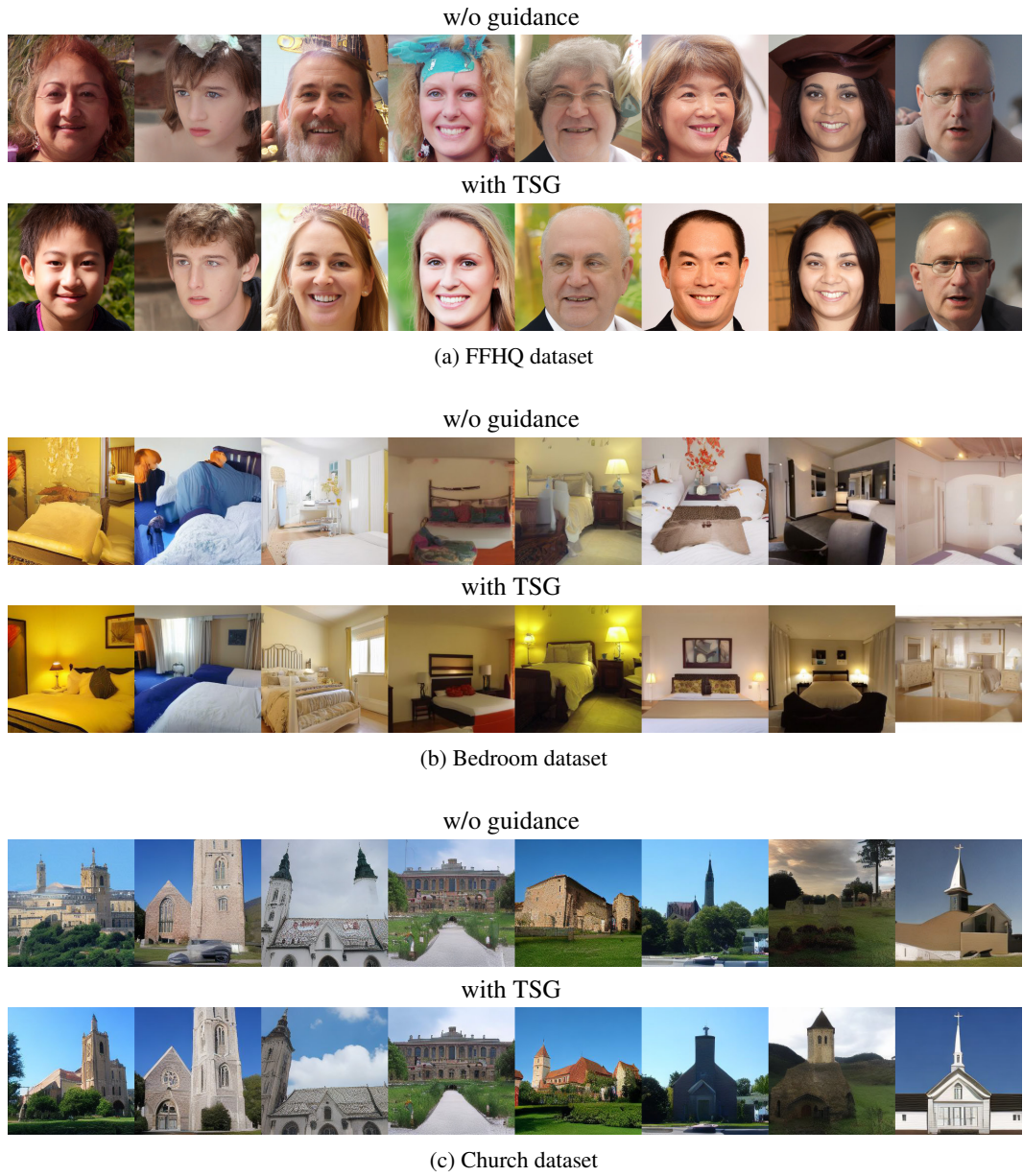


Figure 17: Visual comparison between unguided sampling and sampling with TSG for several unconditional latent diffusion models from Rombach et al. (2022). We observe that TSG consistently improves the quality of all models compared to the baseline sampling.