# Improving Sample Quality of Diffusion Models Using Self-Attention Guidance

Susung Hong        Gyuseong Lee        Wooseok Jang        Seungryong Kim

Korea University, Seoul, Korea

{susung1999, jpl358, jws1997, seungryong_kim}@korea.ac.kr

(a) ADM [7] without (top) and with (bottom) SAG        (b) Stable Diffusion [31] without (top) and with (bottom) SAG

Figure 1: **Qualitative comparisons between unguided (top) and self-attention-guided (bottom) samples.** Unlike classifier guidance (CG) [7] or classifier-free guidance (CFG) [16], self-attention guidance (SAG) does not necessarily require an external condition, *e.g.*, a class label or text prompt, nor additional training, improving the details of the images generated by pre-trained diffusion models such as (a) unconditional ADM [7] and (b) Stable Diffusion [31] with an empty prompt.

## Abstract

*Denoising diffusion models (DDMs) have attracted attention for their exceptional generation quality and diversity. This success is largely attributed to the use of class- or text-conditional diffusion guidance methods, such as classifier and classifier-free guidance. In this paper, we present a more comprehensive perspective that goes beyond the traditional guidance methods. From this generalized perspective, we introduce novel condition- and training-free strategies to enhance the quality of generated images. As a simple solution, blur guidance improves the suitability of intermediate samples for their fine-scale information and structures, enabling diffusion models to generate higher quality samples with a moderate guidance scale. Improving upon this, Self-Attention Guidance (SAG) uses the intermediate self-attention maps of diffusion models to enhance their stability and efficacy. Specifically, SAG adversarially blurs only the regions that diffusion models attend to at each iteration and guides them accordingly. Our experimental re-sults show that our SAG improves the performance of various diffusion models, including ADM, IDDPM, Stable Diffusion, and DiT. Moreover, combining SAG with conventional guidance methods leads to further improvement.*

## 1. Introduction

Recently, denoising diffusion models (DDMs) [35, 37, 14, 7, 15, 31], which synthesize images from noise through an iterative denoising process, have been actively researched and attracted attention due to their exceptional performance in synthesizing high-quality and diverse images.

Behind this remarkable success lies the introduction of diffusion guidance methods [7, 23, 16]. Several studies have revealed that to improve the quality of image samples generated by diffusion models, guidance techniques using

---

The project page and code can be accessed at:
https://ku-cvlab.github.io/Self-Attention-Guidance/

class labels [7, 16] or captions [23] are essential. However, despite the significant improvement provided by these guidance methods, they are bounded within the limits of using external conditions. For example, classifier guidance (CG) [7] requires the training of an additional classifier, and classifier-free guidance (CFG) [16] adds complexity to the training process through label-dropping. In addition, both methods are limited by their need for hard-earned external conditions, which binds them to conditional settings.

In light of the limitations mentioned above, in this work, we present a more general formulation of diffusion guidance that can make use of information within the intermediate samples of diffusion models. This formulation detaches the necessary condition of traditional approaches [16, 7, 23], *i.e.*, the requirement for external information, from diffusion guidance, and facilitates a flexible and condition-free approach to guide diffusion models. This broadens the applicability of diffusion guidance to cases with or without external conditions.

Based on the generalized formulation and the intuition that any internal information within intermediate samples can also serve as guidance, we firstly propose blur guidance as a straightforward solution to improve sample quality. Blur guidance uses the eliminated information resulting from Gaussian blur to guide intermediate samples, exploiting the benign property of Gaussian blur that it naturally removes fine-scale details [17, 20, 30]. While our results show that this method improves sample quality with a moderate guidance scale, it becomes problematic with a large guidance scale, since it may introduce structural ambiguity in entire regions, which makes it difficult to align the prediction of the degraded input with that of the original one.

To improve the effectiveness and stability of blur guidance with a larger guidance scale, we explore the self-attention mechanism of diffusion models. Generally, recent diffusion models [14, 7, 24, 31, 15, 27] are equipped with a self-attention module [40, 8] within their architecture. Claiming that the self-attention is a key to capture salient information during generation process [18, 45, 46, 12], we present Self-Attention Guidance (SAG), which adversarially blurs the region that contains salient information using the self-attention map of diffusion models and guides diffusion models with the residual information. Leveraging the attention maps during the reverse process of diffusion models, it can encouragingly boost the quality and reduce the artifacts through self-conditioning without requiring external information nor additional training, as shown in Fig. 1. The pseudocode and pipeline are provided in Alg. 1 and Fig. 2(b), respectively.

In experiments, we evaluate the effectiveness of the proposed approach by plugging it into various diffusion models including ADM [7], IDDPM [24], Stable Diffusion [31], and DiT [27], which demonstrates our method's broad applicability. We also show that in addition to the increased sample quality when using SAG alone, performance further improves when using it on top of existing guidance schemes, *i.e.*, classifier [7] or classifier-free [16] guidance, demonstrating the orthogonality with the existing methods. Finally, we present ablation studies to validate our choices.

To sum up, our work has the following contributions:

- Generalizing conditional guidance methods [7, 16, 23] into a condition-free method that can be applied to any diffusion model without external conditions, expanding the applicability of guidance.

- Introducing novel guidance, dubbed Self-Attention Guidance (SAG), that uses the internal self-attention maps of diffusion models, improving sample quality without external conditions or additional fine-tuning.

- Demonstrating the orthogonality of SAG to existing conditional models and methods, enabling its flexible combination with others to achieve higher performance.

- Presenting extensive ablation studies to justify the design choices and demonstrate the effectiveness of the proposed method.

## 2. Related Work

**Denoising diffusion models.** Diffusion models [35], which are closely related to score-based models [37, 38], have attracted much attention owing to their superior sampling quality and diversity. As a pioneering work, DDPM [14] generates an image through an iterative process that progressively performs denoising to recover an image. Following this work, there have been several approaches to improve the sampling process, in terms of quality and speed [36, 24, 31, 15, 7]. Notably, IDDPM [24] additionally predicts the variance of the reverse process of the diffusion model. DDIM [36] accelerates the sampling speed by introducing the non-Markovian diffusion process. LDM [31] reduces the computational cost by processing the diffusion process in the latent space.

**Sampling guidance for diffusion models.** Recent works have proposed diffusion guidance methods based on class labels to generate images with higher quality [7, 16]. Classifier guidance (CG) [7] is an approach that uses a trained classifier that guides the reverse process toward a specific class distribution. As an alternative strategy without an additional classifier, Ho and Salimans [16] propose classifier-free guidance (CFG). Due to its simplicity of implementation and effectiveness, the guidance has been used in various high-quality diffusion models [29, 31, 39, 41, 23, 33]. Adopting the concepts of the guidance methods above, Nichol *et al.* [23] propose text-to-image generation with CLIP [28] guidance and CFG. However, these approaches

**Algorithm 1** Self-Attention Guidance (SAG) Sampling

**Functions**:

Model($\mathbf{x}_t$): a diffusion model that outputs the predicted noise $\epsilon_t$, variance $\Sigma_t$, and self-attention map $A_t$ given the input $\mathbf{x}_t$.

Gaussian-Blur($\hat{\mathbf{x}}_0$): a Gaussian blurring function.

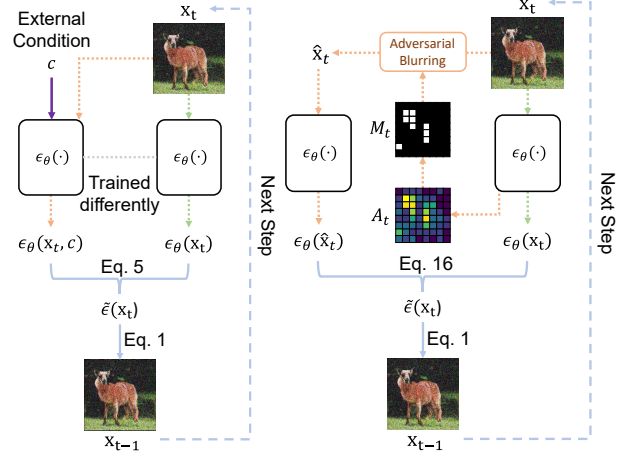$\quad \mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$
$\quad$ **for** $t$ in $T, T-1, ..., 1$ **do**
$\quad\quad \epsilon_t, \Sigma_t, A_t \leftarrow$ Model($\mathbf{x}_t$)
$\quad\quad M_t \leftarrow \mathbb{1}(A_t > \psi)$
$\quad\quad \hat{\mathbf{x}}_0 \leftarrow (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_t)/\sqrt{\bar{\alpha}_t}$ $\quad$ // Eq. 2
$\quad\quad \tilde{\mathbf{x}}_0 \leftarrow$ Gaussian-Blur($\hat{\mathbf{x}}_0$)
$\quad\quad \tilde{\mathbf{x}}_t \leftarrow \sqrt{\bar{\alpha}_t}\tilde{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t$
$\quad\quad \hat{\mathbf{x}}_t \leftarrow (1 - M_t) \odot \mathbf{x}_t + M_t \odot \tilde{\mathbf{x}}_t$ $\quad$ // Eq. 15
$\quad\quad \hat{\epsilon}_t \leftarrow$ Model($\hat{\mathbf{x}}_t$)
$\quad\quad \tilde{\epsilon}_t \leftarrow \hat{\epsilon}_t + (1 + s)(\epsilon_t - \hat{\epsilon}_t)$ $\quad$ // Eq. 16
$\quad\quad \mathbf{x}_{t-1} \sim \mathcal{N}(\frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\tilde{\epsilon}_t), \Sigma_t)$ $\quad$ // Eq. 1
$\quad$ **end for**
$\quad$ **return** $\mathbf{x}_0$



(a) Classifier-free guidance $\quad$ (b) Self-attention guidance

Figure 2: **Comparison of classifier-free guidance [16] and self-attention guidance (SAG).** Compared to classifier-free guidance that uses external class information, SAG extracts the internal information with the self-attention to guide the models, making it training- and condition-free.

have limitations since they do not apply to unlabeled datasets and require additional training procedures [7, 16].

**Self-attention in generative models.** A self-attention mechanism is the key ingredient of Transformer-based models [40]. Notably, it has become a *de facto* method in natural language processing tasks [40] for its expressive power and capability to encode global context, which has inspired many works to incorporate this mechanism into computer vision [8, 18, 45, 46]. Among those, Jiang *et al.* [18] and Zhang *et al.* [45, 46] attempt to bring self-attention into generative adversarial networks (GANs) for better image quality. Following this, diffusion models have also brought self-attention into their model architectures. DDPM [14] initiates this trend by introducing a self-attention layer at a coarse resolution of the U-Net [32]. Inspired by this work, Dhariwal and Nichol [7] measure the boost performance according to the varying number of self-attention heads and resolutions. Concurrently, DiT [27] even accomplishes high performance leveraging Transformer-based backbones.

**Internal representations of diffusion models.** Motivated by the success of diffusion models in generation tasks, some works have tried to utilize the representations of diffusion models to do other tasks, such as semantic segmentation. Brempong *et al.* [2] show that the denoising pre-training boosts the performance on semantic segmentation, and Baranchuk *et al.* [1] propose a label-efficient strategy for semantic segmentation using the U-Net [32] representations of diffusion models. While specific tasks such as text-driven manipulation using cross-attention has been researched concurrently [12], these are inherently different from improving and self-conditioning general diffusion models in a condition-free way leveraging the internal self-attention maps, which is mainly discussed in this paper.

## 3. Preliminaries

**Denoising diffusion probabilistic models.** DDPM [14] is a model that recovers an image from white noise through an iterative denoising process. Formally, given an image $\mathbf{x}_0$ and a variance schedule $\beta_t$ at a timestep $t \in \{T, T-1, \ldots, 1\}$, we can obtain $\mathbf{x}_t$ through the forward process which is defined as a Markovian process. Similarly, given a trained diffusion model parameterized by $\epsilon_\theta(\mathbf{x}_t, t)$ and $\Sigma_\theta(\mathbf{x}_t, t)$, we can define the reverse process. In this case, we set $\Sigma_\theta(\mathbf{x}_t, t)$ to $\sigma_t^2 = \beta_t$ [14] although it is possible to predict the variance [24, 7]. Specifically, given $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$ and $\Sigma_\theta(\mathbf{x}_t, t)$, DDPM samples $\mathbf{x}_{T-1}, \mathbf{x}_{T-2}, \ldots, \mathbf{x}_0$ by computing:

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(\mathbf{x}_t, t)) + \sigma_t \mathbf{z}, \quad (1)$$

where $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i$, $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$, and $\epsilon_\theta$ denotes a neural network parameterized by $\theta$. Note that for simplicity, we define $\epsilon_\theta(\mathbf{x}_t) := \epsilon_\theta(\mathbf{x}_t, t)$ for the rest of the paper. Using the reparameterization trick, we can obtain $\hat{\mathbf{x}}_0$, an intermediate reconstruction of $\mathbf{x}_0$ at a timestep $t$, using the following equation:

$$\hat{\mathbf{x}}_0 = (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_\theta(\mathbf{x}_t, t))/\sqrt{\bar{\alpha}_t}. \quad (2)$$

**Classifier guidance and classifier-free guidance.** To bring the capability of GANs' trading diversity for fidelity to diffusion models, Dhariwal and Nichol [7] propose the classifier guidance that uses an additional classifier $p(c|\mathbf{x}_t)$, where $c$ is a class label. The guidance can be formulated as

the following with a guidance scale $s > 0$:

$$\tilde{\epsilon}(\mathbf{x}_t, c) = \epsilon_\theta(\mathbf{x}_t, c) - s\sigma_t\nabla_{\mathbf{x}_t}\log p(c|\mathbf{x}_t), \quad (3)$$

where $\epsilon_\theta(\mathbf{x}_t, c)$ is a conditional diffusion model, and $\tilde{\epsilon}(\mathbf{x}_t, c)$ is the guided output by the classifier. On the other hand, Ho and Salimans [16] present a classifier-free guidance strategy that achieves the similar effect as classifier guidance without the use of an additional classifier:

$$\tilde{\epsilon}(\mathbf{x}_t, c) = \epsilon_\theta(\mathbf{x}_t, c) + s(\epsilon_\theta(\mathbf{x}_t, c) - \epsilon_\theta(\mathbf{x}_t)) \quad (4)$$

$$= \epsilon_\theta(\mathbf{x}_t) + (1 + s)(\epsilon_\theta(\mathbf{x}_t, c) - \epsilon_\theta(\mathbf{x}_t)). \quad (5)$$

However, this method still demands hard-earned labels and confines the application to conditional diffusion models that use external conditions such as class or text [23, 31, 33] conditions. Moreover, it requires additional training detail that occasionally zero-outs the class embedding in the training phase [16], thus imposing extra complexity.

**Self-attention in diffusion models.** Several works of diffusion models use the U-Net structure [32] with self-attention [40] at one or some of the intermediate layers [14, 7]. Moreover, very recently, diffusion models using Transformers [40] as the backbone has also been proposed [27]. Specifically, for the height $H$ and width $W$, given any feature map $X_t \in \mathbb{R}^{(HW) \times C}$ at a timestep $t$, the $N$-head self-attention is defined as:

$$Q_t^h = X_t W_Q^h, \quad K_t^h = X_t W_K^h, \quad (6)$$

$$A_t^h = \text{softmax}(Q_t^h (K_t^h)^T / \sqrt{d}), \quad (7)$$

where $W_Q^h, W_K^h \in \mathbb{R}^{C \times d}$ for $h = 0, 1, ..., N - 1$. Each $A_t^h$ is then right multiplied by $V_t^h = X_t W_V^h$, where $W_V^h \in \mathbb{R}^{C \times d}$.

## 4. Generalizing Diffusion Guidance

Although classifier guidance and classifier-free guidance have largely contributed to the conditional generation of diffusion models [7, 16, 23], they depend on external inputs. In this work, we broaden our perspective by extending them to handle both cases: with or without external inputs. We also show how CFG [16] can be integrated into our framework at the end of this section.

At a given timestep $t$, the entire input for a diffusion model comprises a generalized condition represented as $\mathbf{h}_t$, and a perturbed sample $\bar{\mathbf{x}}_t$ that lacks $\mathbf{h}_t$. More specifically, the condition $\mathbf{h}_t$ can encompass internal information within $\mathbf{x}_t$, an external condition, or both. With this definition, the resulting guidance is formulated through the utilization of an imaginary regressor, $p_{\text{im}}(\mathbf{h}_t|\bar{\mathbf{x}}_t)$, which is assumed to predict $\mathbf{h}_t$ given $\bar{\mathbf{x}}_t$. Modifying guidance proposed in prior works [38, 7], we present:

$$\tilde{\epsilon}(\bar{\mathbf{x}}_t, \mathbf{h}_t) = \epsilon_\theta(\bar{\mathbf{x}}_t, \mathbf{h}_t) - s\sigma_t\nabla_{\bar{\mathbf{x}}_t}\log p_{\text{im}}(\mathbf{h}_t|\bar{\mathbf{x}}_t), \quad (8)$$



Figure 3: **Comparison of blur guidance with self-attention guidance (SAG) under a large guidance scale.** Given an extreme guidance scale ($s = 5.0$), blur guidance generates relatively noisy images (top) compared to those generated with SAG (bottom).

where we slightly abuse the notation for $\epsilon_\theta(\bar{\mathbf{x}}_t, \mathbf{h}_t)$ since we assume that the inputs are simply aggregated to match the original whole input. Intuitively, the gradient of the regressor, $\nabla_{\bar{\mathbf{x}}_t}\log p_{\text{im}}(\mathbf{h}_t|\bar{\mathbf{x}}_t)$, guides generated samples to be more suitable with that information.

With Bayes' rule, $p_{\text{im}}(\mathbf{h}|\bar{\mathbf{x}}_t) \propto p(\bar{\mathbf{x}}_t|\mathbf{h})/p(\bar{\mathbf{x}}_t)$, and the score of an imaginary regressor $p_{\text{im}}(\mathbf{h}_t|\bar{\mathbf{x}}_t)$ is derived:

$$\nabla_{\bar{\mathbf{x}}_t}\log p_{\text{im}}(\mathbf{h}_t|\bar{\mathbf{x}}_t) = -\frac{1}{\sigma_t}(\epsilon^*(\bar{\mathbf{x}}_t, \mathbf{h}_t) - \epsilon^*(\bar{\mathbf{x}}_t)), \quad (9)$$

where $\epsilon^*$ denotes the true score of the regressor. Eventually, this term is plugged into Eq. 8 and produces:

$$\tilde{\epsilon}(\bar{\mathbf{x}}_t, \mathbf{h}_t) = \epsilon_\theta(\bar{\mathbf{x}}_t, \mathbf{h}_t) - s\sigma_t\nabla_{\bar{\mathbf{x}}_t}\log p_{\text{im}}(\mathbf{h}_t|\bar{\mathbf{x}}_t) \quad (10)$$

$$= \epsilon_\theta(\bar{\mathbf{x}}_t, \mathbf{h}_t) + s(\epsilon_\theta(\bar{\mathbf{x}}_t, \mathbf{h}_t) - \epsilon_\theta(\bar{\mathbf{x}}_t)) \quad (11)$$

$$= \epsilon_\theta(\bar{\mathbf{x}}_t) + (1 + s)(\epsilon_\theta(\bar{\mathbf{x}}_t, \mathbf{h}_t) - \epsilon_\theta(\bar{\mathbf{x}}_t)). \quad (12)$$

Note that Eq. 11 induces a constraint that $\bar{\mathbf{x}}_t$ be in-manifold that the diffusion model $\epsilon_\theta$ defines. Also note that CFG [16] is a special case of Eq. 12 where $\bar{\mathbf{x}}_t = \mathbf{x}_t$, $\mathbf{h}_t = c$, and the imaginary regressor $p_{\text{im}}(\mathbf{h}_t|\bar{\mathbf{x}}_t)$ is reduced into the implicit classifier in [16].

Benefiting from this formulation, we can also define diffusion guidance on unconditional models, which have a sole noised image $\mathbf{x}_t$ as an input and no external label [16, 7], by making it self-conditional on visual information within the intermediate samples of the reverse process. In this light, we present comprehensive discussions on how to find appropriate $\mathbf{h}_t$ for unconditional models and according $\bar{\mathbf{x}}_t$, and subsequently propose guidance in Section 5.

## 5. Utilizing the Self-Attention Map to Improve Sample Quality

The derivation presented in Section 4 implies that by extracting salient information $\mathbf{h}_t$ contained in $\mathbf{x}_t$, it is possible to provide guidance to the reverse process of diffusion models. Inspired by this implication, we propose an innovative

| Dataset | Input | # of steps | SAG | FID ($\downarrow$) | sFID ($\downarrow$) | IS ($\uparrow$) | Precision ($\uparrow$) | Recall ($\uparrow$) |
|---------|-------|-----------|-----|------|------|------|-----------|--------|
| ImageNet 256×256 | Uncond. | 250 | ✗ | 26.21 | 6.35 | 39.70 | 0.61 | **0.63** |
|  |  |  | ✓ | **20.08** | **5.77** | **45.56** | **0.68** | 0.59 |
| ImageNet 256×256 | Cond. | 250 | ✗ | 10.94 | 6.02 | 100.98 | 0.69 | **0.63** |
|  |  |  | ✓ | **9.41** | **5.28** | **104.79** | **0.70** | 0.62 |
| LSUN Cat 256×256 | Uncond. | 250 | ✗ | 7.03 | 8.24 | - | **0.60** | **0.53** |
|  |  |  | ✓ | **6.87** | **8.21** | - | **0.60** | 0.50 |
| LSUN Horse 256×256 | Uncond. | 250 | ✗ | 3.45 | 7.55 | - | **0.68** | **0.56** |
|  |  |  | ✓ | **3.43** | **7.51** | - | **0.68** | 0.55 |

Table 1: **50K results of self-attention guidance on ADM [7] pre-trained on 256×256 images.** The best values are in bold.



Figure 4: **High-frequency masks (top) and the self-attention masks (bottom) of the finally generated images.** Note that the frequency masks are calculated after the generation process, while the self-attention masks are accumulated during the entire reverse process.

| Schedule | Objective | Input | SAG | FID ($\downarrow$) |
|----------|-----------|-------|-----|------|
| cosine | $L_{\text{hybrid}}$ | Uncond. | ✗ | 19.2 |
|  |  |  | ✓ | **18.0** |

Table 2: **A 50K result of self-attention guidance on ID-DPM [24] pre-trained on ImageNet 64×64.**
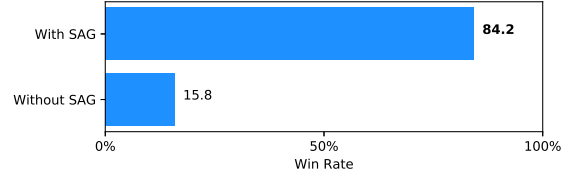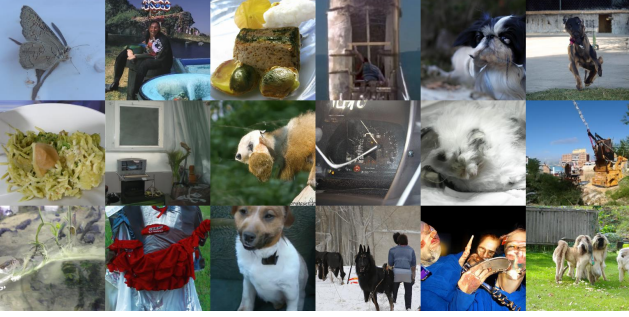


Figure 5: **A result of human evaluation on self-attention guidance with pairs sampled from Stable Diffusion [31].**

guidance technique, Self-Attention Guidance (SAG), which effectively capture the salient information for reverse process while mitigating the risk of out-of-distribution issues of $\bar{\mathbf{x}}_t$ in pre-trained diffusion models. We first explain blur guidance, which is a primitive form of SAG, in Section 5.1, and then we introduce SAG in Section 5.2.

## 5.1. Blur Guidance for Diffusion Models

Gaussian blur is a linear filtering technique that involves convolving an input signal $\hat{\mathbf{x}}_0$ with a Gaussian filter $G_\sigma$ to produce an output $\tilde{\mathbf{x}}_0$. Formally, $\tilde{\mathbf{x}}_0 = \hat{\mathbf{x}}_0 * G_\sigma$, where $*$ represents the convolution operator. As the standard deviation $\sigma$ increases, Gaussian blur reduces the fine-scale details within the input signals and smooths them towards constant [30], resulting in locally indistinguishable ones.

It is evident that there is an information imbalance between $\tilde{\mathbf{x}}_0$ and $\hat{\mathbf{x}}_0$, where $\hat{\mathbf{x}}_0$ contains more fine-scale information. Based on this fundamental insight, we introduce a specialized version of Eq. 12, which we refer to as blur guidance in this paper. In essence, blur guidance intentionally excludes the information from intermediate reconstructions (Eq. 2) during the diffusion process, using this information to guide our predictions towards enhancing the relevance of the images to the information. In detail, blur
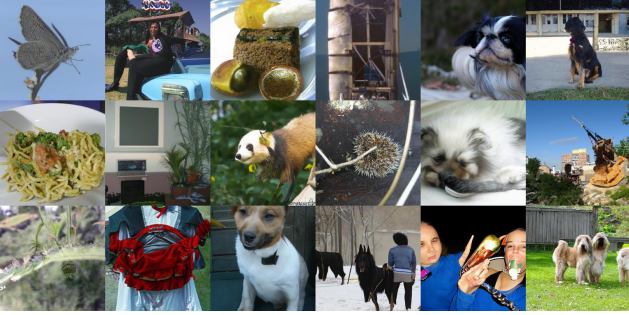
guidance makes the original prediction deviate more from the prediction of the blurred input. Moreover, we note that Gaussian blur has a benign property in that it prevents the resulting signals from deviating significantly from the original manifold with a moderate $\sigma$, i.e., blurring occurs naturally in images [17, 20, 30], which makes Gaussian blur particularly suitable for its application to pre-trained diffusion models. These models generally include latent diffusion models [31], given that the spatial latents also contains low-level information such as local structures [9, 25].

To be specific, we first blur $\hat{\mathbf{x}}_0$ at Eq. 2 with a Gaussian filter $G_\sigma$. Subsequently, we diffuse it again with the noise $\epsilon_\theta(\mathbf{x}_t)$ to produce $\tilde{\mathbf{x}}_t$. It is important to note that by doing this, we bypass the side effect of blur that reduces Gaussian noise, making the guidance depend on the intermediate content rather than the random noise. For brevity and to incorporate diffusion models in latent space [31], we let $\mathbf{x}_t$ represent either noised images or the spatial latents [9, 31].

The blur guidance is then incorporated into Eq. 12 by setting $\bar{\mathbf{x}}_t = \tilde{\mathbf{x}}_t$ and $\mathbf{h}_t = \mathbf{x}_t - \tilde{\mathbf{x}}_t$. In practice, the joint input $(\tilde{\mathbf{x}}_t, \mathbf{h}_t)$ is simply computed as the summation $\mathbf{x}_t = \tilde{\mathbf{x}}_t + \mathbf{h}_t$. The term $\mathbf{x}_t - \tilde{\mathbf{x}}_t$ retains the information present before the blurring process, thus guiding the diffusion process to be more appropriate to the removed salient information in the

(a) Results from ADM [7].


(b) Results from ADM [7] with our method (SAG).

Figure 6: **Uncurated samples from ADM [7] without and with our method (SAG).** Both results are sampled from unconditional ADM pre-trained on ImageNet 256×256 [6], and share the same random seed. The samples guided by SAG typically show fewer artifacts, benefiting from the self-conditioning of the internal conditions.

original input. Our results, as shown in Table 5 "Global", demonstrate the effect of blur guidance in improving the baseline in terms of the quality metrics.
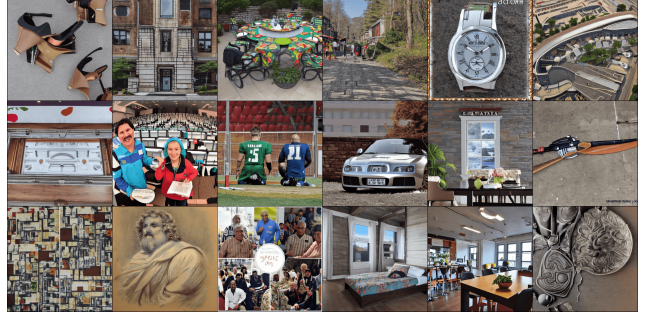
Despite its benefit with moderate guidance scales, the application of blur guidance on existing models with large guidance scales ($s > 5.0$) produces noisy results, as shown in the top row of Fig. 3. We assume that this is because global blur introduces structural ambiguity across entire regions. This makes it difficult to align the prediction of the degraded input with that of the original, contributing to the noisy outcome accumulated over $t$. This issue highlights the need for a more adaptive approach that can capture finer and more relevant information during the reverse process than the global blurring.

## 5.2. Self-Attention Guidance for Diffusion Models

The self-attention mechanism [8, 40] has been shown to be a key component of diffusion models [7, 14]. This mechanism, which is implemented in the backbones of diffusion models, allows the model to attend to salient parts of the input during the generative process [18, 45, 46, 12]. A particular example of the information capture is illustrated in Fig. 4, which shows that the region of the self-attention masks from ADM [7] overlaps with the high-frequency de-


(a) Results from Stable Diffusion [31].


(b) Results from Stable Diffusion [31] with our method (SAG).

Figure 7: **Uncurated samples from Stable Diffusion [31] without and with our method (SAG).** Both results are sampled from Stable Diffusion, and share the same random seed. The prompt is set to an empty prompt with a single space (" ").

tails that diffusion models ought to elaborate on and that are one of essential factors of image generation [4, 42, 43] and human perception [5]. See the appendix for more examples and analyses.

Building upon this intuition, we propose SAG, which leverages the self-attention maps of diffusion models. In essence, we adversarially blur self-attended patches of $\mathbf{x}_t$, *i.e.*, conceal the information of patches that diffusion models attend to. We then use the concealed information to guide diffusion models. In addition, it can be shown that $\bar{\mathbf{x}}_t$ of self-attention guidance contains intact regions of $\mathbf{x}_t$, which means that it does not cause the structural ambiguity of the inputs and thus mitigates the problem of global blur.

To obtain the aggregated self-attention map from Eq. 7, we conduct global average pooling (GAP) to aggregate the stacked self-attention maps $A_t^S \in \mathbb{R}^{N \times (HW) \times (HW)}$ to the dimension $\mathbb{R}^{HW}$, followed by reshaping to $\mathbb{R}^{H \times W}$ and subsequent nearest-neighbor upsampling to match the resolution of $\mathbf{x}_t$:

$$A_t = \text{Upsample}(\text{Reshape}(\text{GAP}(A_t^S))). \quad (13)$$

Generalizing blur guidance, given a masking threshold $\psi$, which is practically set to the mean value of $A_t$, SAG blurs only the masked patches of $\mathbf{x}_t$ according to the self-

"A girl showing a smiling face."  "A living area with a television and a table."  "A Scottish Fold playing with a ball."  " "

Figure 8: **Text-to-image results of Stable Diffusion [31], where the top row is sampled with only CFG and the bottom row with CFG and SAG.** SAG helps the model generate a high-quality image that is more self-conditioned and has fewer artifacts even with an empty prompt (4th column), exhibiting independence from external information.

attention map and is formulated as follows:

$$M_t = \mathbb{1}(A_t > \psi), \tag{14}$$

$$\widehat{\mathbf{x}}_t = (1 - M_t) \odot \mathbf{x}_t + M_t \odot \tilde{\mathbf{x}}_t, \tag{15}$$

$$\tilde{\epsilon}(\mathbf{x}_t) = \epsilon_\theta(\widehat{\mathbf{x}}_t) + (1 + s)(\epsilon_\theta(\mathbf{x}_t) - \epsilon_\theta(\widehat{\mathbf{x}}_t)), \tag{16}$$

where $\odot$ denotes the Hadamard product and $\tilde{\mathbf{x}}_t$ is obtained in the same manner as that in Sec. 5.1. Note that Eq. 16 is also a special case of Eq. 12 where $\mathbf{h}_t = M_t \odot \mathbf{x}_t - M_t \odot \tilde{\mathbf{x}}_t$, $\bar{\mathbf{x}}_t = \widehat{\mathbf{x}}_t$, and the joint input undergoes the simple summation as in Sec. 5.1. Unlike blur guidance, $\widehat{\mathbf{x}}_t$ explicitly contains intact patches of $\mathbf{x}_t$, preventing the output $\epsilon_\theta(\widehat{\mathbf{x}}_t)$ from deviating too far from the original with even a large scale (Fig. 3) as well as effectively concealing the information critical for the reverse process in an adversarial manner.

## 6. Experiments

### 6.1. Experimental Settings

For the experiments, we use two servers with 8 NVIDIA GeForce RTX 3090 GPUs each to sample from. We build upon the pre-trained models of ADM [7], IDDPM [24], Stable Diffusion [31], and DiT [27]. We take all the weights for our experiments from their publicly available repositories, and use the same evaluation metrics as [7], including FID [13], sFID [22], IS [34], and Improved Precision and Recall [19].

### 6.2. Experimental Results

**Unconditional generation with SAG.** We show the effectiveness of SAG on the unconditional models, which demonstrates our condition-free property that CG and CFG do not possess. We use unconditionally pre-trained

| CG [7] | SAG | FID ($\downarrow$) | sFID ($\downarrow$) | Precision ($\uparrow$) | Recall ($\uparrow$) |
|---|---|---|---|---|---|
| ✗ | ✗ | 5.91 | 5.09 | 0.70 | **0.65** |
| ✓ | ✗ | 2.97 | 5.09 | 0.78 | 0.59 |
| ✗ | ✓ | 5.11 | **4.09** | 0.72 | **0.65** |
| ✓ | ✓ | **2.58** | 4.35 | **0.79** | 0.59 |

Table 3: **Compatibility of SAG with CG [7].** The results are from ADM trained on ImageNet 128×128.

| Model | CFG [16] | SAG | FID ($\downarrow$) |
|---|---|---|---|
| DiT-XL/2 [27] | ✓ | ✗ | 2.27 |
|  | ✓ | ✓ | **2.16** |

Table 4: **Compatibility of SAG with CFG [16].** The results are from DiT-XL/2 trained on ImageNet 256×256.

| Masking strategy | FID ($\downarrow$) | IS ($\uparrow$) |
|---|---|---|
| Baseline | 5.98 | 141.72 |
| Global (blur guidance in Sec. 5.1) | 5.82 | 143.15 |
| High-frequency | 5.74 | 148.87 |
| Random | 5.68 | 148.99 |
| Square | 5.68 | 146.50 |
| Self-attention (SAG in Sec. 5.2) | **5.47** | **151.12** |
| DINO [3]-attention | 5.63 | 146.18 |

Table 5: **Ablation study of the masking strategy.** The results are from ADM trained on ImageNet 128×128.

ADM [7] and IDDPM [24] for this experiment, and evaluate 50k samples for the metrics.

We evaluate pre-trained ADM [7] on ImageNet [6] 256×256, LSUN Cat [44], and LSUN Horse [44]. As shown in Table 1, we observe that SAG consistently improves the FID, sFID and IS of unconditional, while it lowers the recall. As explained in recent studies [7, 16], we suspect for the lower recall that there also exists a trade-off relationship between sample fidelity and diversity. Nevertheless, the qualitative improvement is made due to the self-conditioning of our method, as we can see the comparison of unselected samples in Fig. 6.

Subsequently, we include the results of the unconditional model of IDDPM [24] equipped with the proposed method, which is trained on ImageNet at resolution 64×64. The result is in Table 2, which also shows an improvement in terms of FID by applying SAG.

**Conditional generation with SAG.** While our method is effective on unconditional models, Eq. 12 implies the condition-agnosticity, meaning that SAG can also be applied to conditional models. To evaluate SAG on conditional models, we perform an experiment on ADM [7] that is conditionally trained on ImageNet 256x256. The results are presented in Table 1, which demonstrates a similar effect on conditional models as on unconditional ones.

**Stable Diffusion with SAG.** We compare our results with Stable Diffusion [31] using human evaluation (see the ap-

(a) FID  (b) sFID

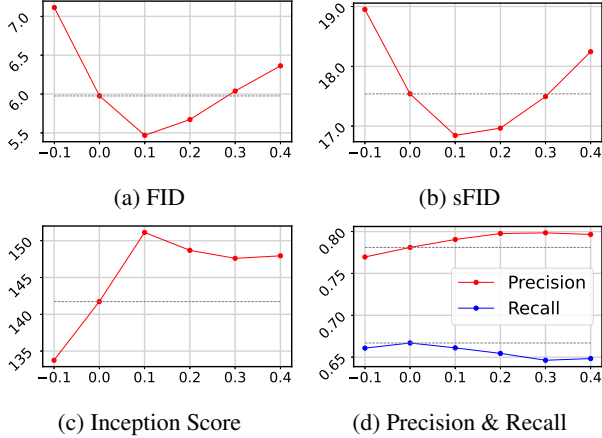(c) Inception Score  (d) Precision & Recall

Figure 9: **Ablation study of the guidance scale.** The x-axis is guidance scale, and the dotted line denotes the performance of the baseline, *i.e.*, the model without SAG. The results are from ADM trained on ImageNet 128×128.

pendix for the protocol) on 500 pairs of images with and without SAG. We use empty prompt for Stable Diffusion and the same random seed for each pair. The results show samples with SAG are more visually favorable or realistic to human. See Fig. 5 and Fig. 7.

In addition, we also broaden the range to text-to-image (T2I) generation by utilizing Stable Diffusion and fusing CFG with SAG, although SAG is not intended for the specific T2I task. Notably, in Fig. 8, the image samples generated from the model with SAG show higher quality and fewer artifacts due to the self-conditioning effect of SAG. Interestingly, even with an empty prompt (Fig. 7 and Fig. 8 4th column), we observe an obviously improved quality. This corroborates the independence of SAG with an external condition.

## 6.3. Ablation Studies and Analyses

**Orthogonality with CG and CFG.** Although designing SAG for unguided models, we can combine it with CG [7] that utilizes external conditions to further improve the performance. To this end, we test four cases to use the guidance, with or without CG and SAG. The metrics are evaluated on 50k samples generated by the ImageNet 128×128 model [7]. As shown in Table 3, we observe additional improvements in FID and precision when using both of them, yet in terms of sFID only giving SAG is the best. This implies that SAG have an orthogonal component with and can be used simultaneously with traditional guidance.

Moreover, CFG [16] is another method of providing class-conditional guidance. However, it requires diffusion models to be trained in a specific manner. Therefore, we use DiT-XL/2 [27], a Transformer [40]-based model which has self-attention layers as well. The 50k results are presented in Table 4. They show that samples guided by CFG

| $\sigma$ | Baseline ($\sigma \to 0$) | $\sigma = 1$ | $\sigma = 3$ | $\sigma = 9$ | $\sigma = 27$ | Avg. pixel ($\sigma \to \infty$) |
|---|---|---|---|---|---|---|
| FID ($\downarrow$) | 5.98 | 5.58 | **5.47** | 5.70 | 5.80 | 5.84 |
| IS ($\uparrow$) | 141.72 | 145.85 | **151.12** | 148.70 | 147.83 | 147.52 |

Table 6: **Ablation study of the sigma ($\sigma$) of Gaussian blur.** The results are from ADM trained on ImageNet 128×128.

| | No guidance | SAG | CFG [16] |
|---|---|---|---|
| GPU memory | 12,167MB | 12,209MB | 12,218MB |
| Run-time | 108.27s | 186.60s | 190.27s |

Table 7: **Computational cost.**

also benefit from the self-conditioning effect of SAG. Note that the combined effect of SAG and CFG is also corroborated by text-to-image samples in Fig. 8.

**Masking strategy.** We test various masking strategies to verify the effectiveness of our self-attention masking with 10k samples on ADM [7]. Those strategies replace the masking function of SAG at each timestep. For a fair comparison, we mask 40% of the pixels of the image for the other masking schemes, which is the equivalent portion of the masked area when the threshold of the self-attention masking is 1.0. The results are in Table 5. We find that the self-attention masking strategy outperforms other masking strategies. Notably, applying global masking, *i.e.* blur guidance, shows the worst performance among the schemes, which validates the motivation for SAG. In addition, we applied the high-frequency mask using FFT on $\hat{x}_0$, as well as the self-attention mask of DINO [3]. However, these methods demonstrated worse performance than ours in terms of FID and IS metrics. Therefore, this result indicates that the self-attention masking is a sufficiently effective method.

**Guidance scale.** We also evaluate the performance changes as the guidance scale changes with 10k samples on ADM [7]. As shown in Fig. 9, we test the scales of $-0.1$, $0.1$, $0.2$, $0.3$, and $0.4$ to ADM and obtain the best FID, sFID, and Inception Score at the guidance scale $s = 0.1$. The precision metric shows the best results when the guidance scale is $s = 0.3$. We also find out that applying self-attention guidance with a negative scale ($s = -0.1$) or a scale that is too large ($s \geq 0.4$) harms the sample quality.

**Gaussian blur.** We examine the effect of changes on $\sigma$ using 10k samples, testing for $\sigma \in \{1, 3, 9, 27\}$ and the extreme cases. As $\sigma \to \infty$, the filter gradually blurs the signal content, reducing every pixel to the average value. Conversely, if $\sigma \to 0$, the signal remains unchanged. The results are in Table 6. SAG is robust against linear changes in $\sigma$, while there still exists an optimal $\sigma$ that yields the best performance. Note that the impact also depends on the input resolution; for instance, a higher input resolution generally requires a larger $\sigma$.

**Computational cost.** We report the computational cost of SAG and CFG [16] in Table 7. The memory and time consumption of SAG is almost the same as CFG, which indicates that the overhead due to the operations in SAG (*e.g.*, blurring and masking) is negligible. However, due to the additional step, the cost is high compared to no guidance.

## 7. Conclusion

We present a novel and general formulation of guidance that utilizes internal information within diffusion models for synthesizing high-quality images. Our method, self-attention guidance, is condition- and training-free, and can be applied to various diffusion models, such as ADM, IDDPM, Stable Diffusion, and DiT, improving their quality and reducing the artifacts via self-conditioning. The results of our experiments demonstrate the effectiveness of our proposed method and the orthogonality of self-attention guidance to existing guidance methods. With the findings and the generalization of guidance, we believe that our work opens new avenues for further research in the field of denoising diffusion models and their guidance.

## Acknowledgements

## References

[1] Dmitry Baranchuk, Andrey Voynov, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Label-efficient semantic segmentation with diffusion models. In *ICLR*, 2021. 3

[2] Emmanuel Asiedu Brempong, Simon Kornblith, Ting Chen, Niki Parmar, Matthias Minderer, and Mohammad Norouzi. Denoising pretraining for semantic segmentation. In *CVPR*, pages 4175–4186, 2022. 3

[3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pages 9650–9660, 2021. 7, 8, 13

[4] Yuanqi Chen, Ge Li, Cece Jin, Shan Liu, and Thomas Li. Ssd-gan: measuring the realness in the spatial and spectral domains. In *AAAI*, volume 35, pages 1105–1112, 2021. 6

[5] Kanjar De and V Masilamani. Image sharpness measure for blurred images in frequency domain. *Procedia Engineering*, 64:149–158, 2013. 6

[6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. Ieee, 2009. 6, 7, 19

[7] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*, 34:8780–8794, 2021. 1, 2, 3, 4, 5, 6, 7, 8, 11, 12, 13, 14, 16, 19, 20, 21

[8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. 2, 3, 6

[9] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, pages 12873–12883, 2021. 5

[10] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *CVPR*, pages 10696–10706, 2022. 15

[11] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages 2961–2969, 2017. 14

[12] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 2, 3, 6

[13] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 30, 2017. 7

[14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020. 1, 2, 3, 4, 6, 11, 12

[15] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1–33, 2022. 1, 2

[16] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. 1, 2, 3, 4, 7, 8, 9, 12, 15

[17] Emiel Hoogeboom and Tim Salimans. Blurring diffusion models. *arXiv preprint arXiv:2209.05557*, 2022. 2, 5

[18] Yifan Jiang, Shiyu Chang, and Zhangyang Wang. Transgan: Two pure transformers can make one strong gan, and that can scale up. *NeurIPS*, 34:14745–14758, 2021. 2, 3, 6

[19] Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *NeurIPS*, 32, 2019. 7

[20] Sangyun Lee, Hyungjin Chung, Jaehyeon Kim, and Jong Chul Ye. Progressive deblurring of diffusion models for coarse-to-fine image synthesis. In *NeurIPS 2022 Workshop on Score-Based Methods*, 2022. 2, 5

[21] Chenlin Meng, Ruiqi Gao, Diederik P Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. *arXiv preprint arXiv:2210.03142*, 2022. 15

[22] Charlie Nash, Jacob Menick, Sander Dieleman, and Peter Battaglia. Generating images with sparse representations. In *ICML*, pages 7958–7968. PMLR, 2021. 7

[23] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 1, 2, 4

[24] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *ICML*, pages 8162–8171. PMLR, 2021. 2, 3, 5, 7, 12

[25] Bjorn Ommer and Joachim M Buhmann. Learning the compositional nature of visual objects. In *CVPR*, pages 1–8. IEEE, 2007. 5

[26] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 32, 2019. 12

[27] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022. 2, 3, 4, 7, 8, 12

[28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021. 2

[29] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 2

[30] Severi Rissanen, Markus Heinonen, and Arno Solin. Generative modelling with inverse heat dissipation. *arXiv preprint arXiv:2206.13397*, 2022. 2, 5

[31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 1, 2, 4, 5, 6, 7, 12, 15

[32] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 3, 4, 13

[33] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 2, 4

[34] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *NeurIPS*, 29, 2016. 7

[35] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pages 2256–2265. PMLR, 2015. 1, 2

[36] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 2

[37] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *NeurIPS*, 32, 2019. 1, 2

[38] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2020. 2, 4

[39] Zhicong Tang, Shuyang Gu, Jianmin Bao, Dong Chen, and Fang Wen. Improved vector quantized diffusion models. *arXiv preprint arXiv:2205.16007*, 2022. 2, 15

[40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017. 2, 3, 4, 6, 8

[41] Tengfei Wang, Ting Zhang, Bo Zhang, Hao Ouyang, Dong Chen, Qifeng Chen, and Fang Wen. Pretraining is all you need for image-to-image translation. *arXiv preprint arXiv:2205.12952*, 2022. 2

[42] Yiwen Xu, Maurice Pagnucco, and Yang Song. Dhg-gan: Diverse image outpainting via decoupled high frequency semantics. In *ACCV*, pages 3977–3993, 2022. 6

[43] Mengping Yang, Zhe Wang, Ziqiu Chi, and Wenyi Feng. Wavegan: Frequency-aware gan for high-fidelity few-shot image generation. In *ECCV*, pages 1–17. Springer, 2022. 6

[44] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 7, 20, 21

[45] Bowen Zhang, Shuyang Gu, Bo Zhang, Jianmin Bao, Dong Chen, Fang Wen, Yong Wang, and Baining Guo. Styleswin: Transformer-based gan for high-resolution image generation. In *CVPR*, pages 11304–11314, 2022. 2, 3, 6

[46] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *ICML*, pages 7354–7363. PMLR, 2019. 2, 3, 6

# Appendix

In this document, we provide additional details of DDPM [7], implementation details of our method, more analyses and results, and the human evaluation protocol. We also discuss the limitations and future work at the end.

## A. Denoising Diffusion Probabilistic Models

DDPM [14] is a generative model that generates an image from white noise with iterative denoising steps. Given an image $\mathbf{x}_0$ and a variance schedule $\beta_t$ for an arbitrary timestep $t \in \{1, 2, \ldots, T\}$, the forward process of DDPM is defined as a Markov process of the form:

$$q(\mathbf{x}_{t+1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t+1}; \sqrt{1-\beta_t}\mathbf{x}_t, \beta_t\mathbf{I}). \tag{17}$$

Note that we can directly get $\mathbf{x}_t$ from $\mathbf{x}_0$ in the closed form:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I}), \tag{18}$$

where $\alpha_t = 1 - \beta_t$, and $\bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i$. Similarly, the reverse process is defined as:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t)\mathbf{I}), \tag{19}$$

where $\mu_\theta$ and $\Sigma_\theta$ denote neural networks with parameter $\theta$.

For the training phase, with $\Sigma_\theta$ fixed to a constant $\sigma_t^2 = \beta_t$ as in DDPM, $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is compared with the following forward posterior:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_0, \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_0, \mathbf{x}_t), \tilde{\beta}_t\mathbf{I}), \tag{20}$$

where $\tilde{\mu}_t = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}\mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\mathbf{x}_t$, and $\tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$. However, instead of directly comparing $\mu_\theta$ to $\tilde{\mu}_t$, Ho *et al.* [14] discover that it is beneficial to optimize $\epsilon_\theta$ with the following simplified objective after reparameterization:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, \quad \text{where} \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}), \tag{21}$$

$$L_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0, t, \epsilon}[||\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, t)||^2]. \tag{22}$$

For sampling $\mathbf{x}_{t-1} \sim p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$, we can compute the following from $\mathbf{x}_T$ to $\mathbf{x}_0$:

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(\mathbf{x}_t, t)) + \sigma_t\mathbf{z}, \tag{23}$$

where $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$. Rewriting Eq. 21, we can get $\hat{\mathbf{x}}_0$ which is a prediction of $\mathbf{x}_0$ at each timestep with the following formula:

$$\hat{\mathbf{x}}_0 = (\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t}\epsilon_\theta(\mathbf{x}_t, t))/\sqrt{\bar{\alpha}_t}. \tag{24}$$

# B. Additional Implementation Details

## B.1. Environmental setting

For the experiments, we use two servers of 8 NVIDIA GeForce RTX 3090 GPUs each to sample from the pre-trained models of ADM [7], IDDPM [24], Stable Diffusion v1.4 [31], and DiT [27]. We build upon the PyTorch [26] implementation of these models, taking all the weights for our experiments from their publicly available repository.

## B.2. Selective blurring

In practice, we efficiently implement selective blurring in Sec. 5.2. At the first step, we blur the intermediate reconstruction $\hat{\mathbf{x}}_0$ of $\mathbf{x}_t$ [14]. Then, we apply masks $1 - M_t$ and $M_t$ on $\hat{\mathbf{x}}_0$ and the blurred version of $\hat{\mathbf{x}}_0$, respectively. Finally, we aggregate the output and then noise it again with the predicted noise $\epsilon_\theta(\mathbf{x}_t)$ that we use for computing $\hat{\mathbf{x}}_0$ above. This process ends up producing the same $\hat{\mathbf{x}}_t$ as Eq. 15 in the main paper.

## B.3. Combination of SAG and CFG

Naïvely, in order to combine SAG with CFG [16] in Stable Diffusion [31] and DiT [27], we have to compute SAG through the conditional and unconditional models, which requires us four feedforward steps. In practice, the guided prediction of noise can be efficiently calculated as follows:

$$\tilde{\epsilon}(\mathbf{x}_t) = \epsilon_\theta(\mathbf{x}_t, c) + s_{\mathrm{c}}(\epsilon_\theta(\mathbf{x}_t, c) - \epsilon_\theta(\mathbf{x}_t)) + s_{\mathrm{s}}(\epsilon_\theta(\mathbf{x}_t) - \epsilon_\theta(\bar{\mathbf{x}}_t)), \tag{25}$$

where $s_{\mathrm{c}}$ and $s_{\mathrm{s}}$ denote the scales of CFG and SAG, respectively, and $c$ denotes a text prompt.

## B.4. Hyperparameter settings

In Table 8, we report our hyperparameter settings for our experiments. In the ablation studies in the main paper, we set the other parameters to the constants in Table 8, while testing the ablated parameter. Note that $\sigma$ is dependent on the input resolution.

| | Model | Self-attention parameter | | | Gaussian-blur parameter |
|---|---|---|---|---|---|
| | | Guidance scale | Threshold | Layer | $\sigma$ |
| ADM [7] | ImageNet 256×256 (unconditional) | 0.5, 0.8 | 1.0 | Output 2 | 9 |
| | ImageNet 256×256 (conditional) | 0.2 | 1.0 | Output 2 | 9 |
| | LSUN Cat 256×256 | 0.05 | 1.0 | Output 2 | 9 |
| | LSUN Horse 256×256 | 0.01 | 1.0 | Output 2 | 9 |
| | ImageNet 128×128 | 0.1 | 1.0 | Output 8 | 3 |
| IDDPM [24] | ImageNet 64×64 (unconditional) | 0.05 | 1.0 | Output 7 | 1 |
| Stable Diffusion [31] | | 0.75, 1.0 | 1.0 | Middle | 1 |
| DiT [27] | | 0.005 | 1.0 | 13th block | 1 |

Table 8: Hyperparameter settings.

Figure 10: **Comparison between self-attention masks of DINO [3] and ADM [7]:** (a) the self-attention masks extracted from DINO [3], (b) the self-attention masks extracted from ADM [7].
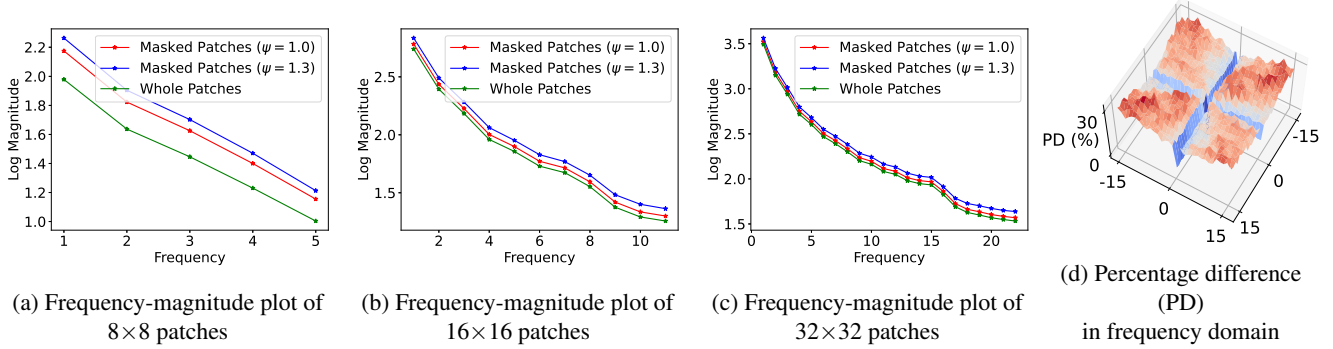


(a) Frequency-magnitude plot of 8×8 patches

(b) Frequency-magnitude plot of 16×16 patches

(c) Frequency-magnitude plot of 32×32 patches

(d) Percentage difference (PD) in frequency domain

Figure 11: **Frequency analysis of the self-attention masks:** (a), (b) and (c) show the frequency-magnitude graphs of 8×8, 16×16, and 32×32 patches, respectively. $\psi$ denotes the masking threshold. (d) is a 3D visualization that shows the percentage difference of magnitude between masked and non-masked patches in the frequency domain regarding the 32×32 patches.

## C. Additional Analyses and Results

### C.1. Exploring the self-attention in diffusion models

We show the visualizations of self-attention maps in the 8×8, 16×16, and 32×32 resolutions of the U-Net [32] of ADM [7] in Fig. 14. The attention maps at $t = 0, 49, 99, 149, 199, 249$ are visualized at each row in order, and the layers are aligned left to right. In this visualization, can see that the attention maps at the intermediate timesteps capture the structure of generated images. Also, we extract the self-attention masks from the different heads and layers from the U-Net and visualize them in Fig. 15 and Fig. 16. *Average* in this figure means the obtained masks after averaging attention maps of the four heads. Moreover, we compare the self-attention masks of ADM with those of DINO [3] in Fig. 10. Compared to the attention masks of DINO, those of ADM are more attending to multiple objects and high-frequency details of the generated images where diffusion models have to elaborate.

Based on the observation, we are interested in two aspects that the self-attention of diffusion models attends to: the frequency and the semantics of the samples. Therefore, we first investigate how the self-attention maps correlate with frequency by comparing the frequency spectra of patches with high attention scores to those of all patches. We observe that high-attention patches contain more high-frequency details (Fig. 11). We then evaluate how the self-attention maps align with foreground objects (Table 9 and Fig. 12) and discover that they capture some semantic information at all resolutions.

### C.2. Additional ablation studies

We conduct experiments on the threshold of self-attention masking that affects the ratio of the blurred region with 10k samples. We test the thresholds of $0.7, 1.0$, and $1.3$. As shown in Table 10, the highest metrics are obtained when the

Figure 12: **Visualization of self-attention masks compared to object masks.** Generated images (top row), the object masks of Mask R-CNN [11] (middle row), and the self-attention masks of unconditional ADM [7] (bottom row).

| Patch size | $\psi$ | Random | Self-attn. | % Diff. |
|---|---|---|---|---|
| 8×8 | 1.0 | 0.16 | 0.23 | + 44% |
|  | 1.3 | 0.09 | 0.14 | + 56% |
| 16×16 | 1.0 | 0.18 | 0.25 | + 39% |
|  | 1.3 | 0.05 | 0.11 | + 120% |
| 32×32 | 1.0 | 0.18 | 0.26 | + 44% |
|  | 1.3 | 0.04 | 0.10 | + 150% |

Table 9: **Semantic analysis of the self-attention masks.** $\psi$ denotes the masking threshold, and % Diff. denotes the percentage difference of the IoU over the random counterpart.

| $\psi$ | Baseline | $\psi = 0.7$ | $\psi = 1.0$ | $\psi = 1.3$ |
|---|---|---|---|---|
| FID ($\downarrow$) | 5.98 | 5.67 | **5.47** | 5.66 |
| IS ($\uparrow$) | 141.72 | 148.60 | **151.12** | 145.58 |

Table 10: **Ablation study of the masking threshold ($\psi$).** The results are derived from ADM trained on ImageNet 128×128.

| Layer | Baseline | In. 11 | In. 8 | Mid. | Out. 2 | Out. 5 | Out. 8 |
|---|---|---|---|---|---|---|---|
| FID ($\downarrow$) | 5.98 | 5.54 | 5.61 | 5.63 | 5.59 | 5.57 | **5.47** |
| IS ($\uparrow$) | 141.72 | 150.07 | 148.20 | 143.44 | 150.62 | 141.73 | **151.12** |

Table 11: **Ablation study of the layer where we extract the attention map.** The results are derived from ADM trained on ImageNet 128×128. We denote the middle block as *Mid.*, and the $n$th layer of the input and output blocks as *In. n* and *Out. n*, respectively.

threshold value is 1.0.

Table 11 shows evaluation results with respect to the attention map extraction layers, evaluated using 10k samples. We select the last self-attention layers of each resolution from the encoder and decoder, and also include the bottleneck layer that divides the encoder and decoder. Regardless of the extraction layer, performance consistently improves over the baseline, while utilizing the self-attention of the final layer yields the best FID and IS results.

### C.3. Qualitative results

In addition to the samples in the main paper, we present random samples with SAG from ADM pre-trained with ImageNet 128×128 (Fig. 17), LSUN Cats (Fig. 18), and LSUN Horse (Fig. 19).

Which row do you think shows the better image quality? 1) The top row 2) The bottom row

Figure 13: **An example of a question.** The participants are not told which row is sampled with our method.

## D. Human Evaluation Protocol

For the human evaluation of SAG with samples from Stable Diffusion [31], we generate 500 pairs with the empty prompt with or without SAG, and the SAG scale is 1.0 for the samples with SAG. Each pair shares the same seed to make it comparable. We show 50 participants 2 groups of 4 samples, one with SAG and the other without SAG, and ask the participants to select a group having higher image quality. An example of a question is in Fig. 13. Neither the pairs are cherry-picked nor filtered. We also do not perform any post-processing with the responses.

## E. Limitations & Future Work

While the increased self-conditioning typically yields results that are more visually appealing to humans, it is important to consider the perspective that the generated images may lack diversity and novelty, a topic that requires discussion. However, at the present stage, the impact of SAG can be effectively moderated by controlling its guidance scale, leading to beneficial applications. Additionally, it requires twice as many feedforward steps, a challenge that is common to CFG [16] and necessitates addressing. A possible solution might involve distilling guidance into diffusion models [21]. This could potentially lessen the computational cost associated with both SAG and CFG, without sacrificing quality.

Moreover, self-attention-based guidance may be more suitable for discrete diffusion models [39, 10], which directly model token probabilities instead of approximating them with continuous values. The integration of these models with our method presents an intriguing topic for future research.
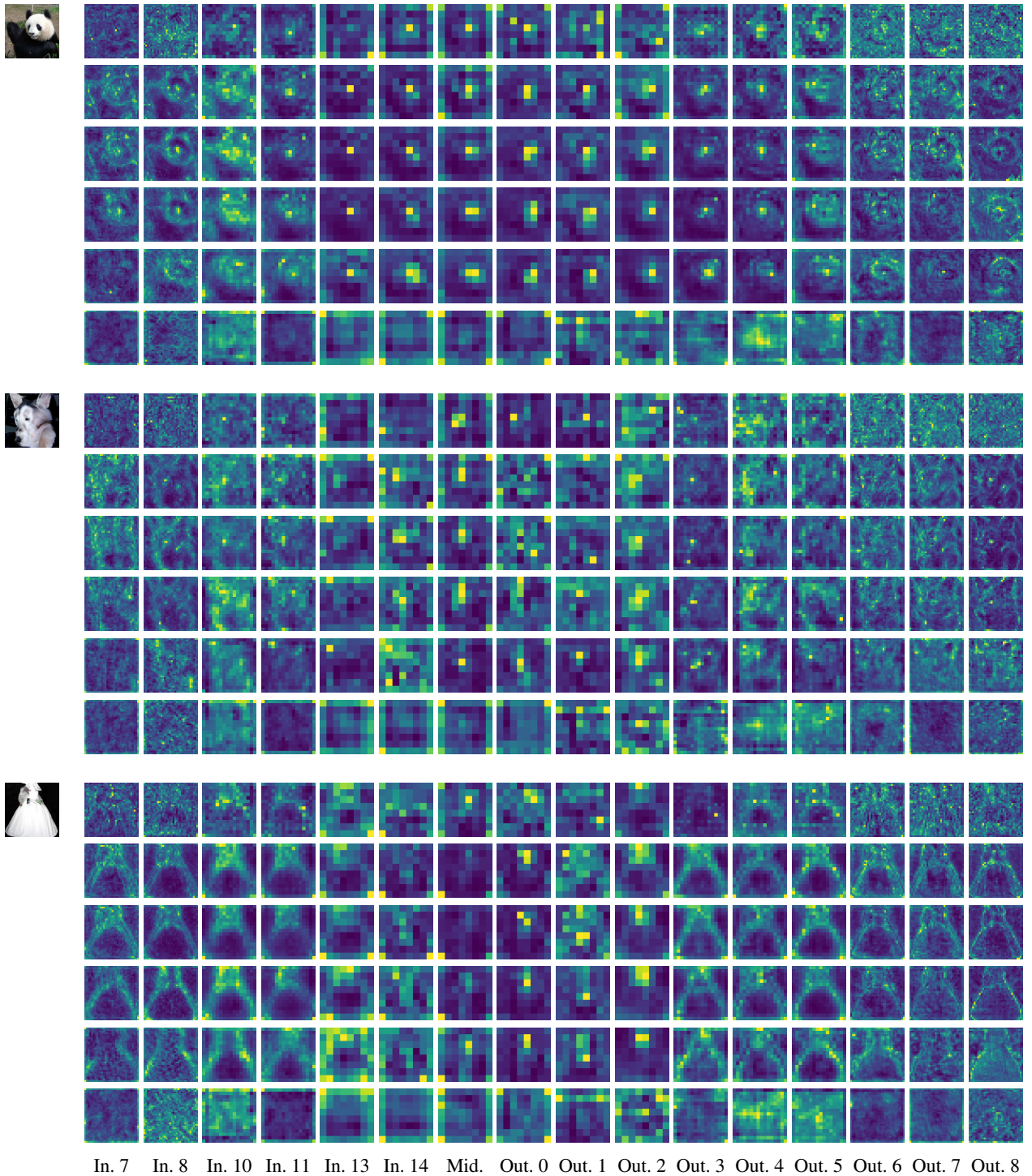
In. 7　In. 8　In. 10　In. 11　In. 13　In. 14　Mid.　Out. 0　Out. 1　Out. 2　Out. 3　Out. 4　Out. 5　Out. 6　Out. 7　Out. 8

Figure 14: **Attention maps at all the self-attention layers of ADM [7].** In. $n$, Mid., and Out. $n$ denote the attention map of the $n$th block of the input blocks, the middle block, and the $n$th block of the output blocks, respectively.
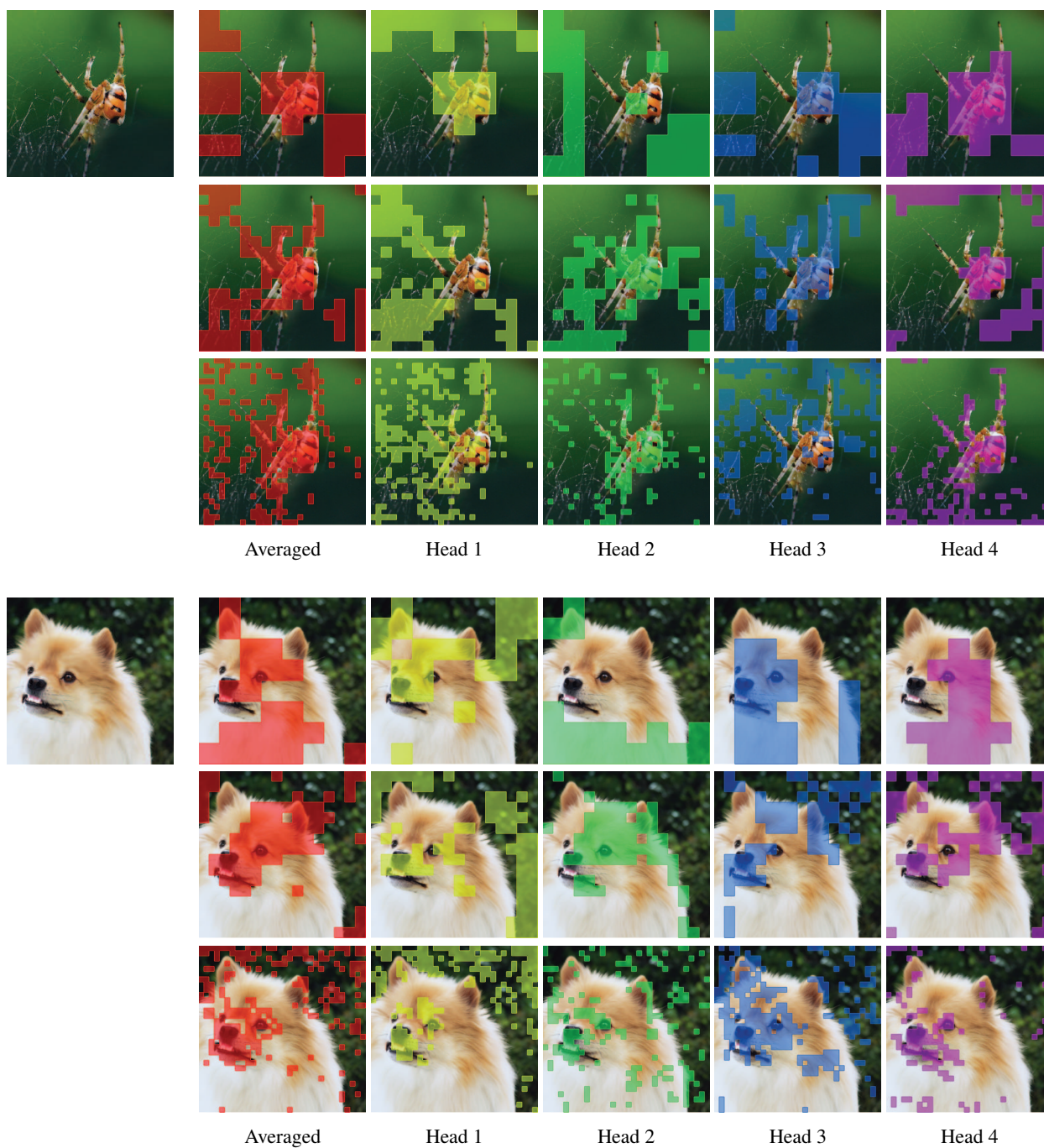
Figure 15: **Visualization of self-attention masks from different layers and heads**. Each row, top to bottom, corresponds to $8 \times 8$, $16 \times 16$ and $32 \times 32$ self-attention layers, respectively.
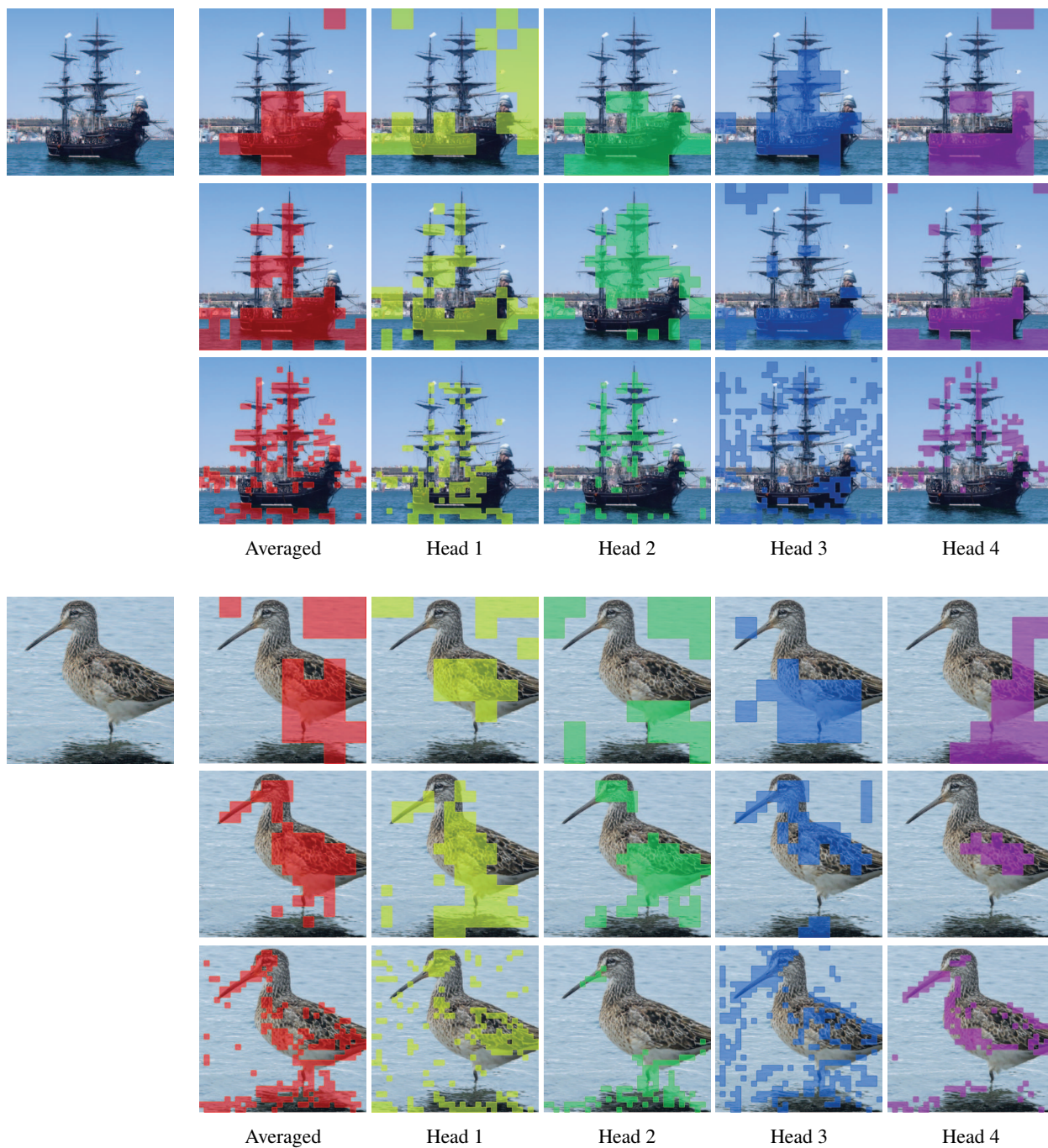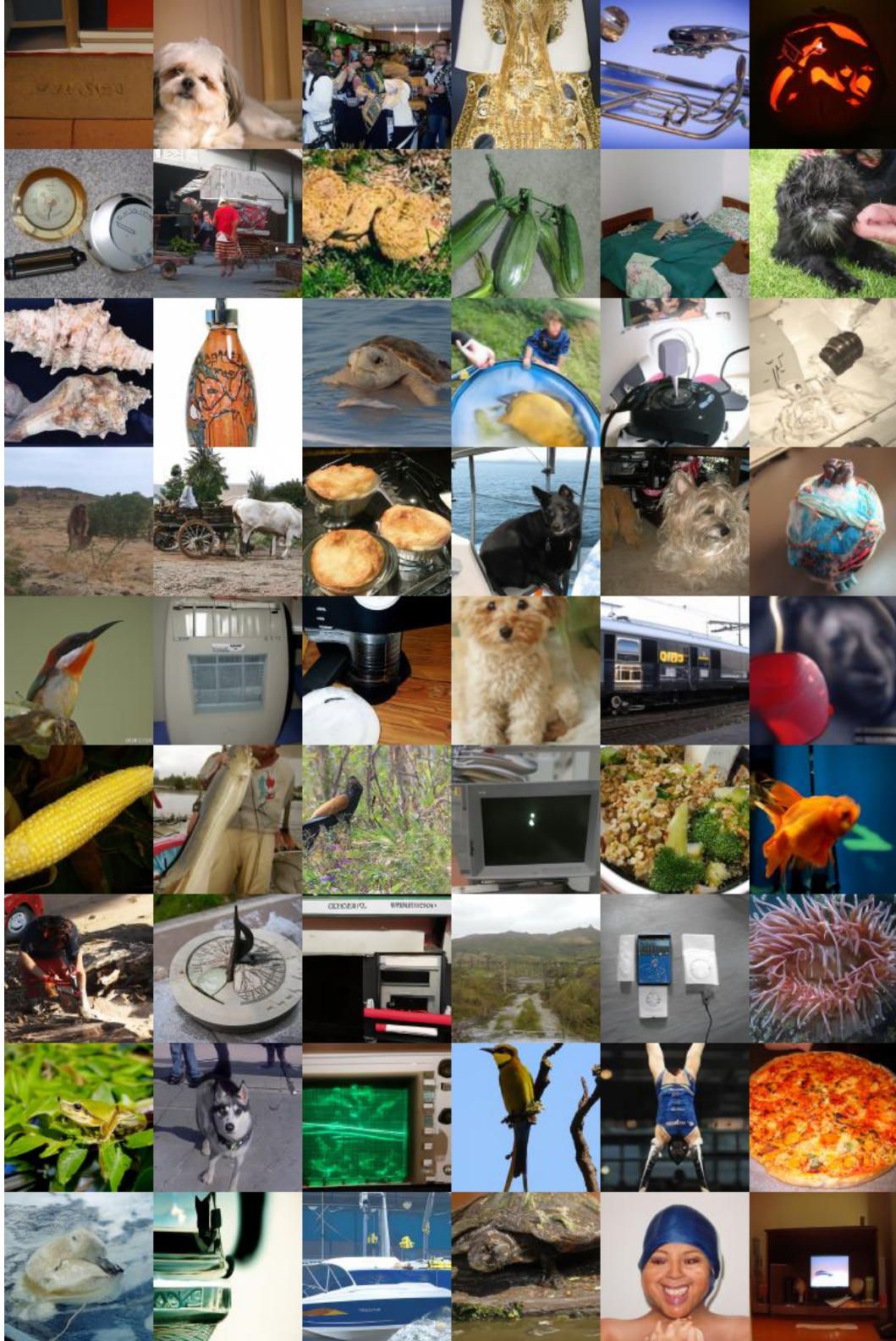
Figure 16: **Visualization of self-attention masks from different layers and heads**. Each row, top to bottom, corresponds to the $8 \times 8$, $16 \times 16$ and $32 \times 32$ self-attention layers, respectively.

Figure 17: **Uncurated samples with our method.** The results are sampled from ADM [7] conditionally pre-trained in ImageNet [6] 128×128 with self-attention and classifier guidance in combination.

Figure 18: **Uncurated samples with our method.** The results are sampled from ADM [7] pre-trained in LSUN Cat [44] with self-attention guidance.
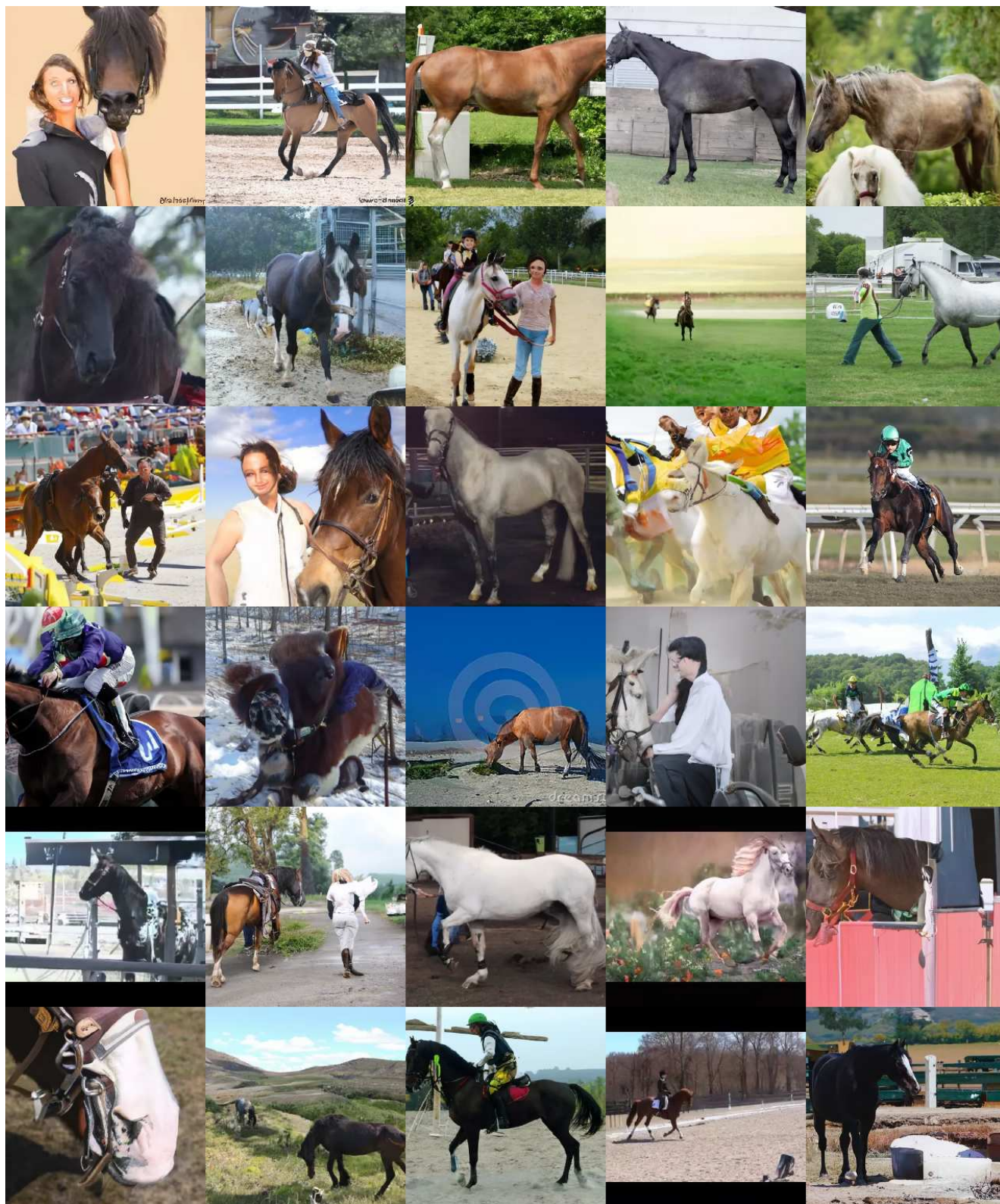
Figure 19: **Uncurated samples with our method.** The results are sampled from ADM [7] pre-trained in LSUN Horse [44] with self-attention guidance.