

Training-Free Layout Control with Cross-Attention Guidance

Minghao Chen Iro Laina Andrea Vedaldi
 Visual Geometry Group, University of Oxford
 {minghao, iro, vedaldi}@robots.ox.ac.uk
[silent-chen.github.io/layout-guidance](https://github.com/silent-chen/layout-guidance)

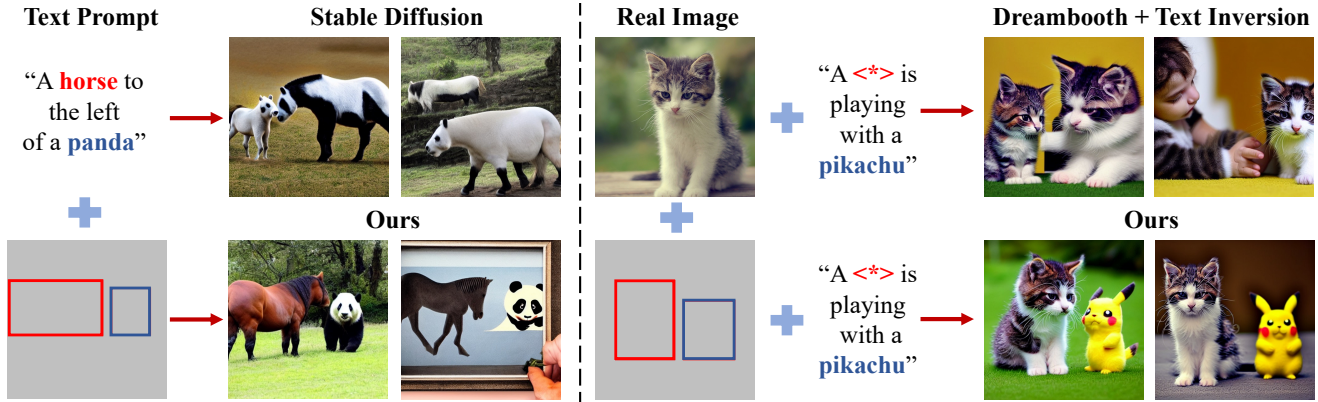


Figure 1. Left: our method controls the layout of an image generated by a pre-trained diffusion model, such as Stable Diffusion [39], without any training or finetuning. It also alleviates the problem of such generators omitting certain objects present in the prompt. Right: given a single real image, our method can be also used to edit the position and context of a subject (represented by $\langle * \rangle$).

Abstract

Recent diffusion-based generators can produce high-quality images from textual prompts. However, they often disregard textual instructions that specify the spatial layout of the composition. We propose a simple approach that achieves robust layout control without the need for training or fine-tuning of the image generator. Our technique manipulates the cross-attention layers that the model uses to interface textual and visual information and steers the generation in the desired direction given, e.g., a user-specified layout. To determine how to best guide attention, we study the role of attention maps and explore two alternative strategies, forward and backward guidance. We thoroughly evaluate our approach on three benchmarks and provide several qualitative examples and a comparative analysis of the two strategies that demonstrate the superiority of backward guidance compared to forward guidance, as well as prior work. We further demonstrate the versatility of layout guidance by extending it to applications such as editing the layout and context of real images.

1. Introduction

Generative AI is one of the most disruptive technologies that emerged in the past years. In computer vision, new text-to-image generation methods, such as DALL-E [37], Imagen [43], and Stable Diffusion [39], have demonstrated that machines are capable of generating images of a quality high enough for use in numerous applications, multiplying the productivity of professional artists as well as lay people.

Despite this success, however, many practical applications of image generation, particularly in a professional setting, require a high level of *control* that such methods lack. Specifications in language-based image generators are textual; and while text can tap into a vast library of high-level concepts, it is a poor vehicle for expressing fine-grained visual nuances in an image. Specifically, text is often inadequate for describing the exact *layout of a composition*.

In fact, as shown in previous work [16], state-of-the-art image generators struggle to correctly interpret simple layout instructions specified via text. For example, when prompting such models with a phrase such as “a dog to the left of a cat”, the “left of” relationship is not always de-

picted accurately in the generated images. In fact, prompts of this nature often cause models to produce erroneous semantics, for example, an image of a cat-dog hybrid. This limitation is exacerbated by unusual compositions, *e.g.*, “horse on top of a house”, which fall outside the typical compositions the model observes during training.

This work provides a better understanding of this limitation and contributes a mechanism to overcome it. To this end, we introduce a method that achieves *layout control* without the need for further training of the image generator, while still maintaining the quality of the generated images.

We note that, while layout cannot be easily controlled via textual prompting, one can *intervene* directly in the cross-attention layers, steering the generation in a direction of choice with user-specified inputs, such as bounding boxes, which we refer to as *layout guidance*. We consider and compare two alternative strategies for such an intervention: “forward guidance” and “backward guidance”. Forward guidance directly biases the cross-attention layers to shift activations in the desired pattern, letting the model incorporate the guidance via the iterated application of its denoising steps. Our main contribution is backward guidance, which uses backpropagation to update the image latents to match the desired layout via energy minimization.

While layout control has already received some attention, with some methods following the forward paradigm [2,45], we show that backward guidance is a more effective mechanism. Our second contribution is then an in-depth investigation of the factors that influence the layout during the image generation process, shedding light on the shortcomings of forward guidance and discussing how backward guidance addresses these. We show that, while there is an intuitive correlation between different concepts and their visual extent, this correlation is more nuanced than one might think, and, perhaps counter-intuitively, even the special tokens in the prompt (start tokens and padding tokens) contribute to shaping the layout.

Finally, we show that our backward guidance outperforms existing methods and seamlessly integrates into applications such as real-image layout editing.

2. Related Work

Text-to-Image Generation. For several years, generative adversarial networks (GANs) [17] have been the dominant approach in image generation from textual prompts [38, 48, 51, 56–58]. Alternative representations for text, such as scene graphs, have also been considered [25]. More recently, the focus has shifted onto text-conditional autoregressive [10, 14, 37, 55] and diffusion models [18, 32, 36, 39, 43], with impressive results in generating images of remarkable fidelity, while avoiding common GAN pitfalls such as training instability and mode collapse [9]. A substantial increase in both the data scale [44] and the size and capabil-

ities of transformer models [35] has played a crucial role in enabling this shift. Typically, these models are designed to accept a textual prompt as input, which may pose a challenge for accurately conveying all details of the image. This problem is exacerbated with longer prompts or when describing atypical scenes. Recent studies have demonstrated the effectiveness of classifier-free guidance [21] in improving the faithfulness of the generations with respect to the input prompt. Others focus on improving compositionality, *e.g.*, by combining multiple diffusion models with different operators [30], and attribute binding [5, 13].

Layout Control in Image Generation. Image generation with spatial conditioning is closely related to layout control and typically done with bounding boxes or semantic maps [12, 33, 46, 47, 52, 60]. These methods do not use text prompts and rely on a closed-set vocabulary to generate images, *i.e.*, the labels of the training distribution (*e.g.*, COCO [29]). Recent image-text models such as CLIP [35] are now enabling the extension to open-vocabulary. However, the precise layout of a composition is still challenging to convey through text alone; even then, the *spatial* fidelity of image generators is extremely limited [16]. Thus, jointly conditioning on text and layout [14, 20, 24] and predicting layout from text [22] have also been considered.

Recent works [1, 2, 4, 6, 27, 45, 50, 53] propose to extend the state-of-the-art Stable Diffusion [39] with spatial conditioning. GLIGEN [27] and ReCo [53] fine-tune the diffusion model with gated self-attention layers and additional regional tokens, respectively. Other works [2, 4, 6, 45, 50] follow a training-free approach. MultiDiffusion [4] adopts the idea from [30] by combining masked noise. eDiff-I [2] and HFG [45] share a similar idea with our forward guidance, directly intervening in the cross-attention. However, they overlook the significance of special tokens in the process. Concurrently with our work, ZestGuide [6] and BoxDiff [50] propose to compute a loss on cross-attention to achieve layout control, which is closer to our backward guidance. Unlike prior work, we use an objective function that does not rely on precise segmentation masks to be provided by the user, and we provide an in-depth analysis of the factors that affect the layout, and consequently, the behavior of both forward and backward strategies. Finally, building on top of diffusion, some recent works show controllable image generation from various other conditioning signals [3, 23, 59], such as depth or edge maps.

Diffusion-Based Image Editing. Most aforementioned methods lack the ability to control or edit an already generated image, or even the ability to edit real images. For example, simply changing a word in the original prompt typically leads to a drastically different generation. This can be circumvented by providing or generating masks for the objects of interest [7, 32]. Prompt-to-prompt [19] addresses this issue with simple text-based edits by exploiting the fact

that the cross-attention layers present in most state-of-the-art architectures connect word tokens to the spatial layout of the generated images. Text-based image editing can also be achieved through single-image model fine-tuning [26, 49]. However, these approaches, while successful at semantically editing entities can only apply such edits *in-place* and do not allow editing of the spatial layout itself.

3. Method

We consider the problem of *layout-guided* text-to-image generation. Text-based image generators allow to sample images $\mathbf{x} \in \mathbb{R}^{3 \times H \times W}$ from a conditional distribution $p(\mathbf{x} | y)$ where y is language description. Given one such generator off-the-shelf, we wish to steer its output to match a desired layout for the generated composition, *without further training or finetuning*. In other words, our objective is to investigate whether pre-trained text-to-image generators can adhere to a layout specified by the user during inference, without having been trained with explicit layout conditioning. In the simplest case, given the text prompt y , the index i of a word y_i in the text prompt, and a bounding box B , we would like to generate an image \mathbf{x} that contains y_i *inside* B , essentially modifying the generator to sample from a new distribution $p(\mathbf{x} | y, B, i)$ with additional controls.

3.1. Preliminaries: Stable Diffusion

We first briefly review the technical details of Stable Diffusion (SD) [39], a publicly accessible, state-of-the-art text-to-image generator representative of an important class of image generators based on diffusion [37, 39, 43]. SD consists of an image encoder and decoder, a text encoder, and a denoising network that operates in latent space.

The text encoder $Y = \phi(y)$ maps the input prompt into a tensor of fixed dimension $Y \in \mathbb{R}^{N \times M}$. This works by prepending a start symbol [SoT] to y and appending $N - |y| - 1$ padding symbols [EoT] at the end, to obtain N symbols in total. Then, the function ϕ , implemented as a large language model (LLM), takes the padded sequence of words as input and produces a corresponding sequence of token vectors $Y_i \in \mathbb{R}^M$ with $i \in \{1, \dots, N\}$ as output.

While not crucial for our discussion, SD’s encoding network h maps images \mathbf{x} to corresponding latent codes $\mathbf{z} = h(\mathbf{x}) \in \mathbb{R}^{4 \times \frac{H}{s} \times \frac{W}{s}}$, where s divides H and W . The function h is an autoencoder with a left inverse h^* , such that $\mathbf{x} = h^* \circ h(\mathbf{x})$. The main purpose of this component is to replace the problem of modeling $p(\mathbf{x} | y)$ with the problem of modeling $p(\mathbf{z} | y)$, reducing the spatial resolution s -fold.

A key component of SD is the iterative conditional denoising network D . This network is trained to output a conditional sample $\mathbf{z} \sim p(\mathbf{z} | y)$ of the latent code \mathbf{z} . It is designed to take a noised sample $\mathbf{z}_t = \alpha_t \mathbf{z} + \sqrt{1 - \alpha_t} \epsilon_t$, as input, where ϵ_t is normally distributed noise and α_t is a decreasing sequence, from $\alpha_0 \approx 1$ to $\alpha_T \approx 0$, representing

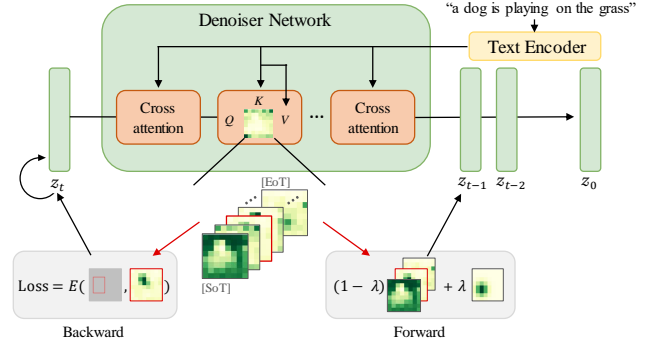


Figure 2. Overview of the two layout guidance strategies. The cross-attention map for a chosen word token is marked with a red border. In forward guidance, the cross-attention maps of the word, start and padding tokens are biased spatially. In backward guidance, we compute instead a loss function and perform backpropagation during the inference process to optimize the latent.

the noise schedule. Then, the network D returns an estimate of the noised sample \mathbf{z}_t : $D(\mathbf{z}_t, y, t) \approx \epsilon_t$. To sample an image, one first samples \mathbf{z}_T , which is normally distributed, and applies D iteratively, to obtain the intermediate codes $\mathbf{z}_{T-1}, \dots, \mathbf{z}_1, \mathbf{z}_0 \approx \mathbf{z}$. Finally, \mathbf{z} is converted back to an image via the image decoder $\mathbf{x} = h^*(\mathbf{z})$.

There is one final aspect of the SD architecture that is relevant for our work. While there are several design choices that make the network D work well in practice, the mechanism that is of interest in our investigation is *cross-attention*, which connects visual and textual information and allows the generation process to be conditioned on text. Each cross-attention layer takes an intermediate feature tensor $\mathbf{z}^{(\gamma)} \in \mathbb{R}^{C \times \frac{H}{r} \times \frac{W}{r}}$ as input, where γ is the index of the relevant layer in the network, and r is a scaling factor defining the spatial resolution at that level of the representation. The cross-attention map $A^{(\gamma)}$ associates each spatial location $u \in \{1, \dots, \frac{H}{r}\} \times \{1, \dots, \frac{W}{r}\}$ to a token indexed by $i \in \{1, \dots, N\}$:

$$A_{ui}^{(\gamma)} = \frac{\exp\langle Q_u^{(\gamma)}, K_i^{(\gamma)} \rangle}{\sum_{j=1}^N \exp\langle Q_u^{(\gamma)}, K_j^{(\gamma)} \rangle}, \quad \mathbf{a}_u^{(\gamma)} = \sum_{i=1}^N A_{ui}^{(\gamma)} V_i^{(\gamma)},$$

where the value $V_i^{(\gamma)}$ and the key $K_i^{(\gamma)}$ are linear transformations of the token embedding Y_i provided by the textual encoder, $Q^{(\gamma)}$ is a linear transformation of $\mathbf{z}^{(\gamma)}$, and $\mathbf{a}_u^{(\gamma)}$ is the output of the cross-attention layer.

3.2. Layout Guidance

Text-to-image generators such as SD struggle to accurately interpret layout instructions provided through text. We thus introduce a method to guide the layout during the generation process by sampling from a distribution $p(\mathbf{x} | y, B, i)$ with additional controls, *e.g.*, user-specified bounding boxes B corresponding to selected text tokens y_i . This

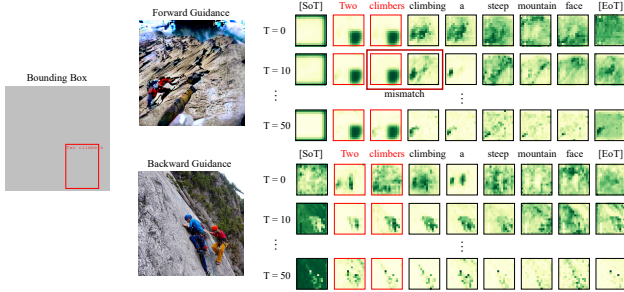


Figure 3. Cross-attention maps during forward and backward guidance. Spatial dependencies between different words negatively affect forward guidance, while backward guidance softly encourages all dependent tokens to match the desired layout.

can be achieved via manipulation of the attention response in certain cross-attention layers in the architecture.

It has already been shown that cross-attention layers regulate the spatial layout of a generated image [19]. Specifically, $A_{ui}^{(\gamma)}$ determines how strongly each location u in layer γ is associated with each of the N text tokens y_i . Since the sum of association strengths $\sum_{i=1}^N A_{ui}^{(\gamma)} = 1$ for each spatial location u , the different tokens can be seen as “competing” for a location. To control the image layout using a bounding box B corresponding to token y_i , the attention can be biased such that locations $u \in B$ within the target box are strongly associated with y_i (while other locations are not). As we discuss below, this can be done without fine-tuning the image generator or training additional layers.

Next, we present a comprehensive investigation of two strategies to achieve training-free layout control: forward and backward guidance (Fig. 2). While instances of *forward* guidance have been discussed in recent work [2, 45], we hereby formalize this approach, identify its limitations, and propose backward guidance as a more effective alternative.

Forward Guidance. In forward guidance, the bounding box B is represented as a smooth windowing function $g_u^{(\gamma)}$ which is equal to a constant $c > 0$ inside the box and quickly falls to zero outside.¹ We rescale the windowing function such that $\|g^{(\gamma)}\|_1 = 1$. Then, we bias a cross-attention map by replacing it with:

$$A_{ui}^{(\gamma)} \leftarrow (1 - \lambda)A_{ui}^{(\gamma)} + \lambda g_u^{(\gamma)} \sum_v A_{vi}^{(\gamma)}, \quad (1)$$

where $\lambda \in [0, 1]$ defines the strength of the intervention. In practice, we normalize the right side of Eq. (1) with a softmax function along the text token dimension, keeping the sum of per-pixel attention equal to 1. Note that (1) only the cross-attention map $A_{:,i}^{(\gamma)}$ of the i -th token is manipulated, and (2) the window is weighed by the mass $\sum_v A_{vi}^{(\gamma)}$ so as to leave the latter unchanged.

¹For simplicity, in our implementation, we put a Gaussian blob with σ decided by the resolution, height, and width of the bounding box.

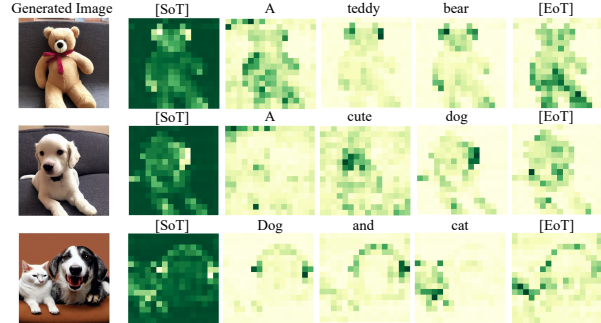


Figure 4. Cross-attention maps of different text prompts at the generation process, indicating that start [SoT] and padding [EoT] tokens carry rich semantic and layout information.

This intervention is applied for a number of iterations of the denoiser network D at selected layers $\gamma \in \Gamma$. This means that the activation maps computed by each selected layer are independently modified following Eq. (1).

A critical analysis reveals that forward guidance is a simplistic approach that suffers from inherent constraints hindering its ability to provide effective layout control. As we discuss in Section 3.3, this is primarily due to various factors that influence the layout during the generation process, including spatial dependencies among text tokens and spatial information “hidden” in the initial noise.

Backward Guidance. To address the shortcomings of forward guidance, we introduce an alternative mechanism, which we refer to as backward guidance. Instead of directly manipulating attention maps, in backward guidance, we bias the attention by introducing an energy function

$$E(A^{(\gamma)}, B, i) = \left(1 - \frac{\sum_{u \in B} A_{ui}^{(\gamma)}}{\sum_u A_{ui}^{(\gamma)}} \right)^2. \quad (2)$$

Optimizing this function encourages the cross-attention map of the i -th token to obtain higher values inside the area specified by B . Specifically, at each application of the denoiser D , when layer $\gamma \in \Gamma$ is evaluated, the gradient of the loss (2) is computed via backpropagation to update the latent $z_t (\equiv z_t^{(0)})$:

$$z_t \leftarrow z_t - \sigma_t^2 \eta \nabla_{z_t} \sum_{\gamma \in \Gamma} E(A^{(\gamma)}, B, i), \quad (3)$$

where $\eta > 0$ is a scale factor controlling the strength of the guidance and $\sigma_t = \sqrt{(1 - \alpha_t)/\alpha_t}$. By updating the latent, the cross-attention maps of all tokens are indirectly influenced by backward guidance. To generate an image, we alternate between gradient updates and denoising steps.

3.3. Analysis and Discussion

Next, we detail a comparative analysis between the forward and backward strategies. To motivate backward guidance and understand its effectiveness, we shed light on the

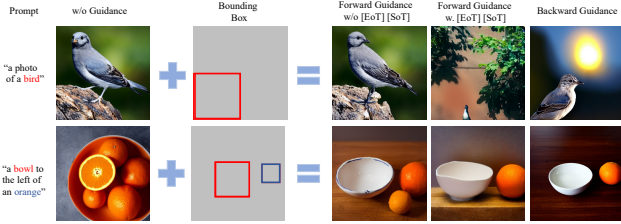


Figure 5. Comparison between forward and backward guidance, including guidance of start and padding tokens.

significance of all tokens and the influence of the initial noise in shaping the layout during the generation process.

The Role of Word Tokens. One important consideration is that the text encoder fuses information from different words when processing a prompt due to self-attention. This results in a “semantic overlap”: information from one token being encoded by another token. In other words, text embeddings capture both word-specific *and contextual* information, *e.g.*, subject-verb-object dependencies. This overlap is then transferred from the text encoder into the diffusion process via the cross-attention layers, resulting in *spatial* overlap. The example in Figure 3 illustrates this overlap in the cross-attention maps of different words. It also shows the behavior of forward and backward guidance when providing spatial conditioning for the phrase “two climbers”. It becomes evident that the mismatch between the attention map of the conditioned phrase and its spatial dependencies with other words (“climbing”, “a”) causes forward guidance to disregard the layout condition. Instead, backward guidance indirectly drives all attention maps toward the layout condition as necessary, because it acts on the latent codes.

The Role of Special Tokens. Another crucial finding is that the cross-attention maps of $[S_{OT}]$ and $[E_{OT}]$ tokens, which do not correspond to content words in the input text, still carry significant semantic and layout information. As we show in Figure 4, the cross-attention maps of $[E_{OT}]$ tokens correspond to salient regions in the generated image, *i.e.*, typically the union of individual semantic entities in the text prompt. $[S_{OT}]$ behaves complementarily to $[E_{OT}]$, emphasizing the background. For forward guidance to be effective, it is thus necessary to intervene not only on selected content tokens but also on the special ones. We use the union of the input boxes as guidance for $[E_{OT}]$ and the reverse for $[S_{OT}]$. However, we have empirically found that this sometimes results in overly aggressive guidance, which harms image fidelity. Backward guidance, on the other hand, does not suffer from such drawbacks, as it optimizes the latent. We discuss this further in the supplement.

The Role of Initial Noise. Finally, the initial noise of the diffusion process plays an important role in shaping the layout of the images. We have empirically observed that the noise contains an intrinsic layout; *e.g.*, when prompting the model with phrases like “an image of a dog” and “an im-

age of a cat” using the same seed, it generates images with consistent layouts, placing the dog and the cat in the same locations. We provide examples in the supplement.

An initial noise with an intrinsic layout close to the one given by users is easier to optimize and results in higher fidelity. Therefore, selecting a noise pattern that aligns with the desired layout can further boost the effectiveness of the guidance. In backward guidance, the loss applied to the cross-attention maps can, in fact, double as a metric for initial noise selection. Specifically, we sample different noise patterns and evaluate Eq. (2) after applying backward guidance for a few steps. This allows us to pick the best-aligned initial noise. Please see the supplement for detailed results.

Forward vs. Backward. In summary, forward and backward guidance use different mechanisms to manipulate cross-attention. Forward guidance *directly* modifies cross-attention to conform to the prescribed pattern, which is “forced” repeatedly for a number of denoising iterations. While it does not incur any extra computational cost, it struggles to provide robust control over the layout, as non-guided tokens may cause the generation to deviate from the desired pattern. In contrast, backward guidance uses a loss function to evaluate whether the attention follows the desired pattern. While slower than forward guidance, backward guidance is more refined, as it indirectly encourages all tokens (guided and non-guided ones) to adhere to the layout through latent updates.

3.4. Real-image Layout Editing

Layout guidance can be used in combination with other techniques that build on diffusion-based image generators. We demonstrate this for the task of real-image editing. To this end, we incorporate backward guidance into two methods that are commonly used for personalization of diffusion models given real images, namely Textual Inversion (TI) [15] and Dreambooth [42]. TI extends an existing image generator with a new concept given one or several images as examples, by optimizing a learnable text token $\langle * \rangle$ for the concept. Dreambooth attempts to capture the appearance of a particular subject of which several images are available by fine-tuning a pre-trained text-to-image model. Then, new images of the learned concept can be generated.

Neither method supports *localized* spatial control over the newly generated images; their edits are usually global and semantic. To achieve this, we apply backward guidance on the Dreambooth-finetuned model and the TI-optimized token as part of a prompt. This allows us to control the layout of the generated images while preserving the identity of the original object represented by $\langle * \rangle$.

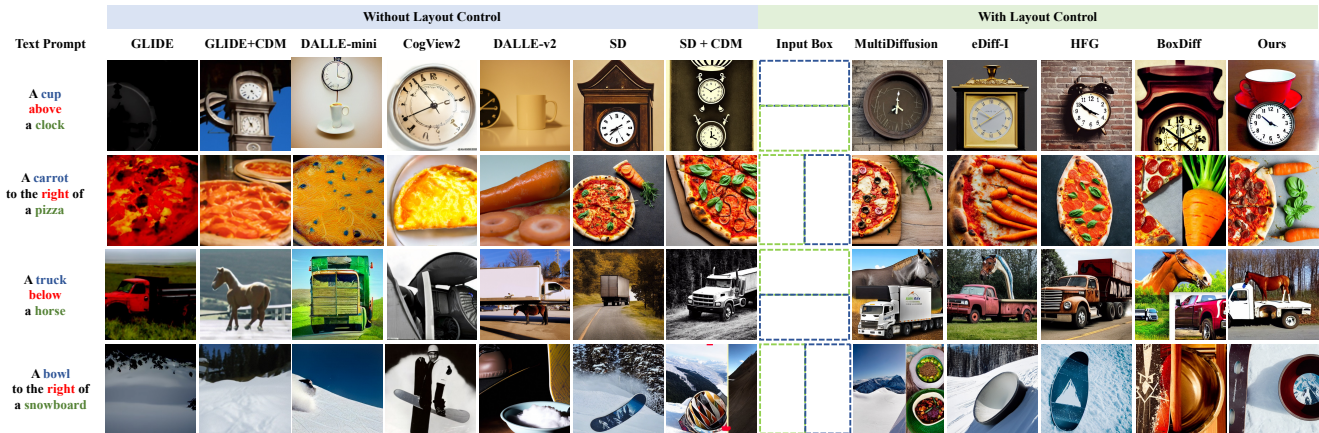


Figure 6. Qualitative comparison of different text-to-image models with text prompts defined in [16]. As stated in [16], current text-to-image models fail to understand spatial relationships without explicit layout conditioning. However, we achieved control of the generated images with the help of guidance on cross-attention maps.

Model	OA (%)	VISOR (%)		Runtime
		uncond	cond	
Stable Diffusion	27.4	16.4	59.8	~ 4 sec/image
Ours (FG)	25.9	23.5	90.7	~ 4 sec/image
Ours (FG*)	27.6	26.1	95.0	~ 5 sec/image
Ours (BG)	38.8	37.6	96.9	~ 8 sec/image
Ours (BG + NS)	43.7	42.3	96.9	~ 9 sec/image

Table 1. Comparison of the forward (FG) and backward (BG) strategies, including noise selection (NS). FG*: forward guidance includes [S_oT] and [E_oT] tokens. We randomly sampled 1000 text prompts and compute metrics based on VISOR [16].

Model	OA (%)	VISOR (%)	
		uncond	cond
GLIDE [32]	3.36	1.98	59.06
GLIDE + CDM [30]	10.17	6.43	63.21
DALLE-mini [8]	27.10	16.17	59.67
CogView2 [11]	18.47	12.17	65.89
DALLE-v2 [36]	63.93	37.89	59.27
SD [39]	29.86	18.81	62.98
SD + CDM [30]	23.27	14.99	64.41
SD + Ours	40.01	38.8	95.95

Table 2. Comparison of backward guidance (ours) with text-to-image generation models based on the VISOR [16] protocol.

4. Experiments

In this section, we evaluate our approach for training-free layout guidance, quantitatively comparing variants of forward and backward guidance and providing comparisons to prior and concurrent work on three benchmarks.

Method	COCO 2014		Flickr30K		
	FID (↓)	mAP (↑)	FID (↓)	AP _P (↑)	mAP (↑)
MultiDiffusion [4]	70.7	22.3	84.1	21.6	11.9
eDiff-I [2]	72.5	21.7	85.3	21.4	9.7
HFG [45]	72.2	21.5	85.6	22.4	10.7
BoxDiff [50]	72.6	24.1	78.7	26.0	16.6
Stable Diffusion [39]	72.3	19.2	76.4	19.4	8.7
Stable Diffusion + Ours	73.3	35.7	78.9	35.6	17.9
GLIGEN [27]	69.1	62.8	77.3	87.2	31.4
GLIGEN + Ours	66.7	65.1	78.1	88.9	32.7

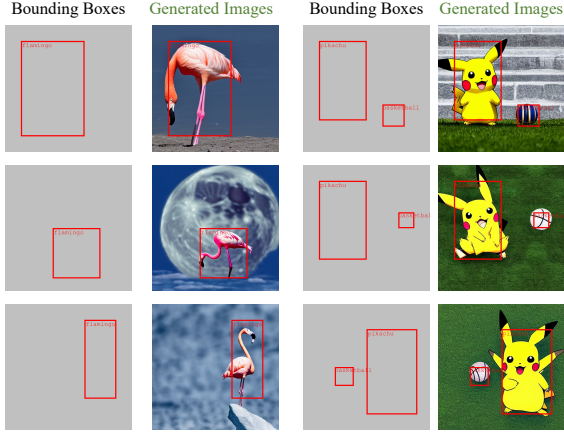
Table 3. Comparison with other layout-to-image models. Our approach improves spatial fidelity (suggested by higher AP/mAP scores). mAP is calculated with an IoU threshold of 0.3.

4.1. Experimental setup

Implementation Details. We utilize Stable-Diffusion (SD) V-1.5 [39] trained on the LAION-5B dataset [44] as the default pre-trained image generator, if not specified. For a detailed description of the architecture and noise scheduler please see the supplement.

For forward guidance, we apply Eq. (1) to every layer of the denoiser network for the first 40 steps of the diffusion process and set $\lambda = 0.8$. For backward guidance, we calculate the loss on the cross-attention maps of the mid-block and the first block of the up-sampling branch of the denoising network (U-Net [41]) as we found this to be the optimal setting to balance control and fidelity. We set $\eta = 30$ by default but found that values between 30-50 work well across most settings. Since the layout of the generated image is typically established in the early stages of inference, backward guidance is performed during the initial 10 steps of the diffusion process and repeated 5 times at each step.

Evaluation Benchmarks. We quantitatively evaluate our approach on three benchmarks: VISOR [16], COCO



A flamingo is standing on the moon. A Pikachu is playing a basketball on grass.

Figure 7. Our method controls the objects inside the generated images with user-specified bounding boxes. On the left, the size and position of *flamingo* changes according to the bounding box. On the right, we show the ability to control multiple objects.

2014 [29], and Flickr30K Entities [34, 54]. We discuss the ethical concerns of the dataset usage in the supp. VISOR proposes metrics to quantify the spatial understanding abilities of text-to-image models. For COCO 2014, we follow the same setup adopted by prior work [4], which uses only a subset of the annotated objects per image. Finally, we introduce the Flickr30K Entities dataset as another benchmark to evaluate layout control, since it contains image-caption pairs with visual grounding. Details for all benchmarks and metrics are provided in the supplementary material.

4.2. Forward vs. Backward Guidance

First, we compare the two different modes of guidance (forward and backward) in Table 1 using the VISOR protocol with 1,000 randomly chosen text samples. The biggest advantage of forward guidance is that the computation overhead is negligible, thus leading to a faster inference time. However, we observe that, compared to (unguided) SD, forward guidance does not significantly increase the object accuracy (OA), while the backward mechanism yields a notably higher OA. In terms of evaluating the generated spatial relationships (VISOR conditional/unconditional metrics), both forward and backward guidance obtain significantly better results than the SD baseline. We also find that the inclusion of [SOT] and [EOT] tokens improves forward guidance, which confirms our analysis and insights in Section 3.3, yet backward guidance still achieves superior performance. Finally, noise selection using backward guidance offers a significant boost in all metrics.

We provide a qualitative comparison of the forward and backward mechanisms in Figure 5, including the impact of special tokens on forward guidance. Backward guidance

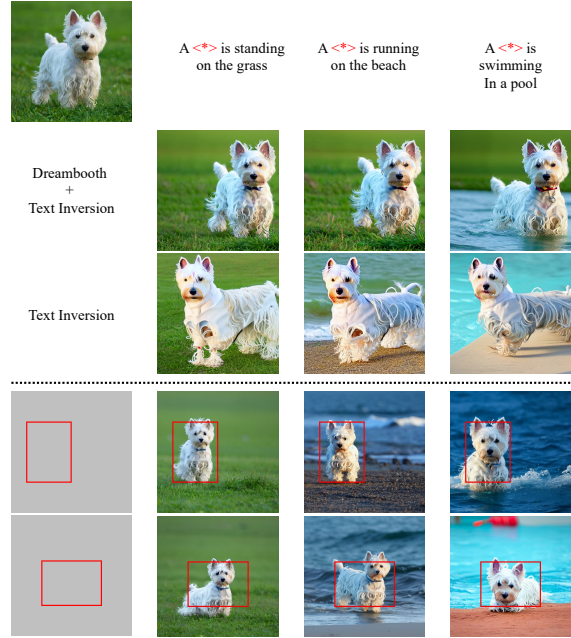


Figure 8. The top left is the real image input. The images above the dash are generations using only text inversion (TI) [15] and Dreambooth [42]. The images under the line are generated by our method on top of Dreambooth and TI.

achieves a better alignment between the generated objects and the input bounding boxes. It also helps to address the issue of objects occasionally being omitted from the generated images in diffusion models.

4.3. Comparisons to Prior Work

In Table 2, we compare our method with text-to-image generation methods that do not use layout control. We note that comparisons are fair since, in this setting (VISOR), manual user input is not required for guidance (see supplement). Our method exhibits remarkable performance under the VISOR_{cond} metric, achieving an accuracy of 95.95%, and higher OA compared to the baseline (SD). Although OA does not directly assess layout, the improvement can be explained by the fact that unguided SD often fails to generate correct semantics in atypical compositions. We also note that, while DALLE-v2 [36] achieves the highest OA overall, it appears to struggle more with layout instructions compared to SD, as indicated by a lower VISOR_{cond} score.

In Table 3, we compare our backward guidance to other mechanisms for layout conditioning. Apart from the entries in the last two rows, all methods are based on Stable Diffusion [40] V1.5. Remarkably, our backward guidance surpasses other layout conditioning methods by a significant margin, achieving over a 9-point improvement in mAP and AP_P on COCO and Flickr30K. Notably, in direct comparison with the concurrent BoxDiff model [50], we achieve gains of 11.6 in mAP and 9.6 in AP_P, all while maintain-

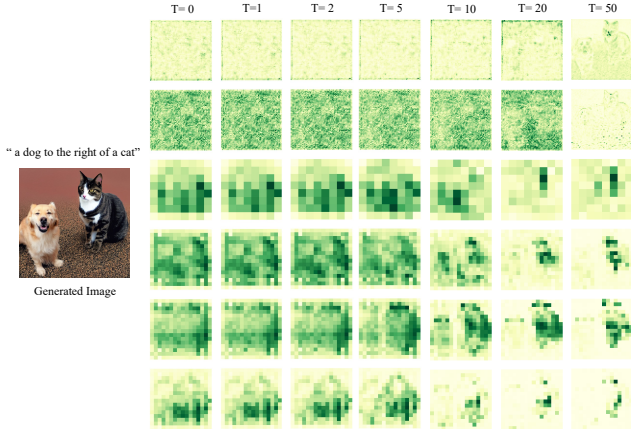


Figure 9. The cross-attention map of the word “cat” at different layers (top to bottom) across different timesteps (left to right).

ing analogous image quality. Finally, we show that our approach can be used complementarily to methods like GLIGEN [27] that train additional layers for layout conditioning, further improving their performance.

In Figure 6, we qualitatively compare different text-to-image models using prompts sampled from [16]. Methods that do not use layout control are not capable of inferring the spatial relationships between objects based purely on textual input and often fail to generate one or both objects. We also observe that even methods with layout conditioning struggle in this setting, especially those that adopt a forward guidance paradigm (eDiff-I [2], HFG [45]). In the case of BoxDiff [50], the lower quality could potentially be due to overlooking the impact of special tokens and the loss function design. In contrast, our approach (backward-guided SD) can accurately position objects within a scene, even when they are rarely seen together, such as “snowboard” and “bowl”, and achieves the best adherence to the prompt without loss of image fidelity. More examples of our approach are shown in Figure 7, demonstrating precise control over the *size* and *position* of one or more objects, including unconventional object categories, such as “flamingo” or “pikachu”, and atypical scene compositions.

4.4. Further Analysis and Applications

Real-Image Layout Editing. We showcase the potential of backward layout guidance for editing real images in Figure 8, confirming its effectiveness at changing the position, gesture, and orientation of the “dog” (based on the aspect ratio of the bounding box) to fit the new context, without altering its identity. As shown in the same figure, the capability to precisely control object size and position cannot be attained through Dreambooth/TI alone. This highlights the potential of our method in a wide range of applications related to image editing and manipulation.

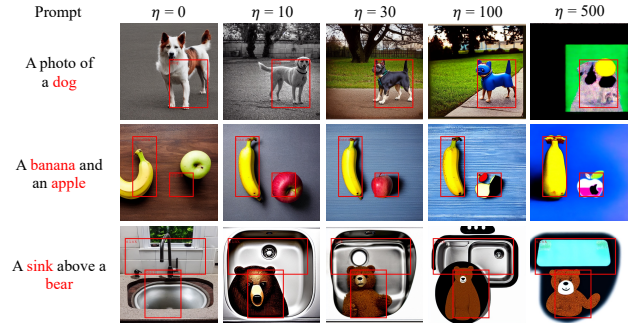


Figure 10. Qualitative comparison of different loss scales in the backward guidance. We increase the loss scale from left to right keeping the same prompt and random seed. With increasing scale, the objects are more tightly constrained inside the bounding boxes. However, for very high scales, fidelity decreases significantly.

Cross-attention Layers and Guidance Steps. We also investigate the layers and the number of guidance steps that are necessary to achieve layout control. Cross-attention maps at various layers of the denoising network are presented in Figure 9. We observe that the first layers of (down-sampling) do not capture much information about the object (here, the “cat”). We found it most effective to perform backward guidance only on the mid and up-sampling blocks of the architecture. The figure also illustrates that object outlines are typically generated in the early steps of the diffusion process, before $T = 20$. Based on our experimentation, we find that 10-20 steps are generally suitable for guidance. Additional quantitative analysis and examples are presented in the supplement.

Loss Scale Factor. In Figure 10, we qualitatively analyze the impact of the loss scale factor η . We observe that increasing the loss weight leads to stronger control over the generated images, but at the cost of some fidelity, particularly with higher scales. The optimal loss scale setting depends on the difficulty of the text prompt. For example, an atypical prompt like “a sink above a bear” requires stronger guidance to generate both objects successfully (without guidance, *i.e.*, $\eta = 0$, the bear is not generated). This suggests that layout guidance helps the generator “recognize” multiple objects in the text prompt.

5. Conclusions

In this paper, we investigated the potential of manipulating the spatial layout of images generated by large, pre-trained text-to-image models without additional training or fine-tuning. Through our exploration, we discovered that both the cross-attention maps and the initial noise of the diffusion play a dominant role in determining the layout and that even the cross-attention maps of special tokens contain valuable semantic and spatial information. We identify and analyze the mechanism behind most prior work: forward

guidance. Moreover, based on our analysis, we propose a new technique “backward guidance” that overcomes the shortcomings of forward guidance. Finally, we demonstrate the versatility of our training-free strategy by extending it to applications such as real-image layout editing.

Ethics. We use the Flickr30K Entities and MS-COCO datasets in a manner compatible with their terms. Some of these images may accidentally contain faces or other personal information, but we do not make use of these images or image regions. For further details on ethics, data protection, and copyright please see <https://www.robots.ox.ac.uk/~vedaldi/research/union/ethics.html>.

Acknowledgements. This research is supported by ERC-CoG UNION 101001212.

References

- [1] Omri Avrahami, Thomas Hayes, Oran Gafni, Sonal Gupta, Yaniv Taigman, Devi Parikh, Dani Lischinski, Ohad Fried, and Xi Yin. Spatext: Spatio-textual representation for controllable image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18370–18380, 2023. 2, 12
- [2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 2, 4, 6, 8
- [3] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. *arXiv preprint arXiv:2302.07121*, 2023. 2
- [4] Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: Fusing diffusion paths for controlled image generation. *arXiv preprint arXiv:2302.08113*, 2023. 2, 6, 7
- [5] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on Graphics (TOG)*, 42(4):1–10, 2023. 2
- [6] Guillaume Couairon, Marlène Careil, Matthieu Cord, Stéphane Lathuilière, and Jakob Verbeek. Zero-shot spatial layout conditioning for text-to-image diffusion models. *arXiv preprint arXiv:2306.13754*, 2023. 2, 12
- [7] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427*, 2022. 2
- [8] Boris Dayma, Suraj Patil, Pedro Cuenca, Khalid Saifullah, Tanishq Abraham, Phúc Le Khac, Luke Melas, and Ritobrata Ghosh. Dall·e mini, 2021. 6
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 2
- [10] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835, 2021. 2
- [11] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. *Advances in Neural Information Processing Systems*, pages 16890–16902, 2022. 6
- [12] Wan-Cyuan Fan, Yen-Chun Chen, DongDong Chen, Yu Cheng, Lu Yuan, and Yu-Chiang Frank Wang. Frido: Feature pyramid diffusion for complex scene image synthesis. *arXiv preprint arXiv:2208.13753*, 2022. 2
- [13] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Reddy Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. In *The Eleventh International Conference on Learning Representations*, 2023. 2
- [14] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pages 89–106. Springer, 2022. 2
- [15] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. 5, 7
- [16] Tejas Gokhale, Hamid Palangi, Besmira Nushi, Vibhav Vineet, Eric Horvitz, Ece Kamar, Chitta Baral, and Yezhou Yang. Benchmarking spatial relationships in text-to-image generation. *arXiv preprint arXiv:2212.10015*, 2022. 1, 2, 6, 8, 12, 19, 20
- [17] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2014. 2
- [18] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10696–10706, 2022. 2
- [19] 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, 4
- [20] Tobias Hinz, Stefan Heinrich, and Stefan Wermter. Generating multiple objects at spatially distinct locations. *arXiv preprint arXiv:1901.00686*, 2019. 2
- [21] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 2, 12
- [22] Seunghoon Hong, Dingdong Yang, Jongwook Choi, and Honglak Lee. Inferring semantic layout for hierarchical text-to-image synthesis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7986–7994, 2018. 2

- [23] Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023. [2](#)
- [24] Xun Huang, Arun Mallya, Ting-Chun Wang, and Ming-Yu Liu. Multimodal conditional image synthesis with product-of-experts gans. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVI*, pages 91–109. Springer, 2022. [2](#)
- [25] Justin Johnson, Agrim Gupta, and Li Fei-Fei. Image generation from scene graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1219–1228, 2018. [2](#)
- [26] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022. [3](#)
- [27] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. *arXiv preprint arXiv:2301.07093*, 2023. [2](#), [6](#), [8](#)
- [28] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014. [14](#)
- [29] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014. [2](#), [7](#), [12](#), [13](#)
- [30] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVII*, pages 423–439. Springer, 2022. [2](#), [6](#)
- [31] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection with vision transformers. *arXiv preprint arXiv:2205.06230*, 2022. [12](#)
- [32] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. [2](#), [6](#)
- [33] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2337–2346, 2019. [2](#)
- [34] Bryan A. Plummer, Liwei Wang, Christopher M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. *International Journal of Computer Vision (IJCV)*, 2017. [7](#), [12](#), [13](#)
- [35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. [2](#)
- [36] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, page 3, 2022. [2](#), [6](#), [7](#)
- [37] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. [1](#), [2](#), [3](#)
- [38] Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *International conference on machine learning*, pages 1060–1069. PMLR, 2016. [2](#)
- [39] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. [1](#), [2](#), [3](#), [6](#), [12](#), [14](#)
- [40] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. [7](#)
- [41] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. [6](#), [12](#)
- [42] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arXiv:2208.12242*, 2022. [5](#), [7](#)
- [43] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. [1](#), [2](#), [3](#)
- [44] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022. [2](#), [6](#)
- [45] Jaskirat Singh, Stephen Gould, and Liang Zheng. High-fidelity guided image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5997–6006, 2023. [2](#), [4](#), [6](#), [8](#)
- [46] Wei Sun and Tianfu Wu. Image synthesis from reconfigurable layout and style. In *Proceedings of the IEEE/CVF*

- International Conference on Computer Vision*, pages 10531–10540, 2019. 2
- [47] Tristan Sylvain, Pengchuan Zhang, Yoshua Bengio, R Devon Hjelm, and Shikhar Sharma. Object-centric image generation from layouts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2647–2655, 2021. 2
- [48] Ming Tao, Hao Tang, Fei Wu, Xiao-Yuan Jing, Bing-Kun Bao, and Changsheng Xu. Df-gan: A simple and effective baseline for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16515–16525, 2022. 2
- [49] Dani Valevski, Matan Kalman, Yossi Matias, and Yaniv Leviathan. Unitune: Text-driven image editing by fine tuning an image generation model on a single image. *arXiv preprint arXiv:2210.09477*, 2022. 3
- [50] Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained diffusion. *arXiv preprint arXiv:2307.10816*, 2023. 2, 6, 7, 8
- [51] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1316–1324, 2018. 2
- [52] Zuopeng Yang, Daqing Liu, Chaoyue Wang, Jie Yang, and Dacheng Tao. Modeling image composition for complex scene generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7764–7773, 2022. 2
- [53] Zhengyuan Yang, Jianfeng Wang, Zhe Gan, Linjie Li, Kevin Lin, Chenfei Wu, Nan Duan, Zicheng Liu, Ce Liu, Michael Zeng, et al. Reco: Region-controlled text-to-image generation. *arXiv preprint arXiv:2211.15518*, 2022. 2
- [54] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *TACL*, 2014. 7, 12
- [55] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Guntjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 2
- [56] Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, and Yinfei Yang. Cross-modal contrastive learning for text-to-image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 833–842, 2021. 2
- [57] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 5907–5915, 2017. 2
- [58] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1947–1962, 2018. 2
- [59] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 2
- [60] Bo Zhao, Lili Meng, Weidong Yin, and Leonid Sigal. Image generation from layout. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8584–8593, 2019. 2
- [61] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022. 12

This appendix contains the following parts:

- **Implementation Details.** We provide more details of the experimental settings, including the network architecture and noise scheduler.
- **Evaluation Dataset and Metrics.** We provide the details of dataset and evaluation metrics used in the experiments part.
- **Ablation Study.** A detailed quantitative evaluation is presented to understand the impact of various components and hyper-parameter selections. We investigate the influence of guided steps, layer-specific losses, and the loss scale factor for backward guidance.
- **Analysis on Initial Noise.** We demonstrate that different prompts with the same initial noise generate images with similar layouts. Therefore, a good choice of initial noise is essential for the success of guidance. Additionally, we quantitatively prove that using the defined loss on cross-attention allows for optimal initial noise selection, enhancing guidance performance.
- **Analysis on Different Tokens.** We visualize the cross-attention map of different prompts and provide extra experiments about controlling the layout of the generated image with only padding tokens.
- **More Examples.** We provide additional examples of our method, including examples under VISOR [16] protocol and real image editing examples.

A. Implementation Details

We provide additional details of our experimental settings.

Network Architecture. In all experiments, we use the Stable Diffusion (SD) V-1.5 [39] as our base model without any architecture modification. The diffusion model is trained in the latent space of an autoencoder. Specifically, the diffusion model adopts the U-Net [41] architecture with a relative downsampling factor of 8. The down-sampling branch of the U-Net has three sequential cross-attention blocks. The mid part of the U-Net has only one cross-attention block. The up-sampling branch of the U-Net has three sequential cross-attention blocks. In each cross-attention block, there are repeated layers following the order: ResBlock \rightarrow Self-Attention \rightarrow Cross-Attention. The cross-attention blocks in the down-sampling branch, mid part, and up-sampling have 2, 1, and 3 such repeated patterns, respectively.

Noise Scheduler. The LMSDscheudler is utilized in all of our experiments with 51 time steps and beta values starting at 0.00085 and ending at 0.012, following a linear scheduler. We also adopt class-free guidance, as suggested

in [21], with a guidance scale of 7.5, consistent with prior work [39].

B. Evaluation Datasets and Metrics

VISOR [16]. We follow the evaluation process described in [16] to compute the VISOR metric, which is designed to quantify the spatial understanding abilities of text-to-image models. This metric focuses on two-dimensional relationships, such as *left*, *right*, *above*, and *below*, between two objects. We measure object accuracy (OA), which is the probability that the generated image contains both objects specified in the text prompt. $VISOR_{uncond}$ is the probability that generating both objects with correct spatial relationship, and $VISOR_{cond}$ is the conditional probability of correct spatial relationships being generated, given that both objects were generated correctly. To generate text prompts for evaluation, we use the 80 object categories from the MS COCO dataset [29], resulting in a total of $80 \times 79 \times 4 = 25,280$ prompts considering any combination of two object categories for each spatial relationship. For each prompt, we generate a single image. As layout guidance inputs we split the image canvas into two, vertically or horizontally, to create two adjacent bounding boxes depending on the type of spatial relationship defined by the text prompt. This only imposes a weak constraint on the layout and can be done automatically (no user intervention is required). For a fair comparison to previous methods that are evaluated in [16], we use the same detection model (OWL-ViT [31]) as in [16] when computing the VISOR metric.

COCO 2014 [29] We randomly sampled 1000 images with their annotations for evaluation from the COCO 2014 validation dataset. The bounding boxes in COCO 2014 are not always grounded in the corresponding caption. Therefore, we append the object labels to the caption as the text prompt for image generation following a similar setting in [1, 6]. Besides, we only pick one to three bounding boxes with areas covering at least 5% of the image panel per sample following the setting in [6]. To assess the quality of the generated images we compute the FID score between the sampled 1000 images from COCO and generated images. We use an open-vocabulary object detector (Detic [61]) to obtain the respective grounding on generated images, which allows quantifying *layout fidelity* using common detection metrics such as average precision (AP). The vocabulary of the detector is constrained to all the COCO object labels.

Flickr30K Entities [34] Finally, we evaluate our method on the Flickr30k Entities dataset [34, 54], which contains image-caption pairs. Since the dataset provides visual groundings of the textual descriptions, we sample a single caption per image and its corresponding bounding boxes and use this as input to perform layout-controlled guidance with SD. We generate a total of 1,000 images using sam-

ples from the validation set. Similarly to the metric used in COCO 2014, we compute the FID score between the original images and the generated ones and use AP as a metric of layout control. To enhance the reliability of the detector, we convert each phrase in the Flickr30 dataset into a single noun (*e.g.*, *ball*) and filter out unrelated nouns, resulting in a total of 303 categories. For each image, the target vocabulary for Detic is defined by the grounded entities in the corresponding caption. To avoid contaminating the evaluation process with perceived human attributes (such as gender, age, occupation, etc.), we also convert all instances of people (man, woman, child, boy, girl, policeman, student, etc.) to the super-class “person” in the target vocabulary for Detic. Since then the *person* category is predominant, we calculate average precision separately for this category (AP_p) but also report the mean average precision across all categories (mAP).

C. Ablation Study

Guidance Step	FID (\downarrow)	AP_p (\uparrow)	mAP (\uparrow)	Inference Time
0	76	19.4	8.7	~ 4sec/image
2	81.2	29.7	13.7	~ 4sec/image
5	81.4	30.3	15.6	~ 6 sec/image
10	82.0	33.5	16.7	~ 8 sec/image
15	82.3	35.5	14.7	~ 10 sec/image
20	83.2	35.6	15.3	~ 12 sec/image
30	83.5	35.7	15.3	~ 15 sec/image

Table A4. Ablation study on guidance steps.

In this section, we supplement the ablation studies in the main paper with quantitative evaluations, studying the impact of the guided steps, loss scale factor, and the effect of backward guidance on different layers of the denoising network. We followed the same setting as described above and in Section 4.1 (main paper) using 1000 captions and their corresponding bounding boxes from the Flickr30K Entities [34] dataset to generate images with a pre-specified layout.

Impact of Guidance Step. Firstly, we explore the effects of guided steps we perform in the diffusion process. The results are shown in Tab. A4, we evaluate image quality (FID), AP_p , layout control (mAP) while varying the number of *guided* steps. We found no improvement in mAP after 10 steps, and FID gradually deteriorates. We hypothesize that this decline may result from potentially shifting the latent vector away from the distribution that corresponds to the original text embedding. Besides, we could see that when increasing the guided steps in the diffusion process, the computation time increases. This is a trade-off question. Generally, a range of 2-10 guidance steps suffices, but users can fine-tune this based on their specific requirements.

Impact of Layers. Secondly, we study the behavior of different layers, by applying backward guidance on the cross-attention maps across different layers of the network. The results are shown in Table A5. As stated in Section 4.4 and illustrated in the table, layers of the down-sampling branch are the least likely to conform to layout control (with $Down-1 < Down-2 < Down-3$ in terms of mAP). In general, high-resolution blocks (such as Down-1 or Up-3) should not be used to control the layout. To achieve the best trade-off between image quality and layout control, a combination of the mid-block (Mid-1) and the first cross-attention block in up-sampling branch (Up-1) of the U-Net is the optimal choice overall.

Impact of Loss Scale Factor. We follow the same setup to evaluate the scale factor η used as the strength of the loss for backward guidance. In Table A6 we report the FID, AP_p and mAP for different loss scale factors. When the loss scale is set to 5–50, the FID is low compared to a larger loss scale factor, indicating that the quality with a loss scale factor of 5–50 is generally good. To achieve better control over the layout, the loss scale factors of 20–50 have the lowest AP_p and mAP. According to the experiments, a loss scale factor of 20–50 works generally well. This factor can be adjusted by the user to get more realistic images or achieve better control over the layout.

D. Analysis on Initial Noise

We conduct an in-depth analysis of the effects of initial noise. As illustrated in Figure A1, the initial noise reveals significant spatial information about the layout. Notably, altering sentence words does not affect this final layout significantly. Figure A2 offers a visual comparison of scenarios with and without noise selection. The results indicate that our backward guidance achieves better control when noise selection is employed. Furthermore, Table A7 quantitatively assesses the impact of noise selection on COCO 2014 and Flickr30K datasets. Methods incorporating noise selection consistently outperform others, underscoring the efficacy of our loss as a noise selection metric.

E. Analysis on Different Tokens

Next, we study the type of information carried by different tokens and their corresponding cross-attention maps, which is relevant for layout guidance.

Removing Word Tokens. We first show that the *padding* tokens convey a significant amount of semantic information. In Figure A3, we randomly pick a subset of captions from MSCOCO [29] and generate images using the Stable Diffusion model and the full caption as the input prompt. As a comparison, after the captions pass through the text encoder, we replace the token embeddings of each caption with the embeddings of its corresponding padding tokens,

Base Model	Down-1	Down-2	Down-3	Mid-1	Up-1	Up-2	Up-3	FID (↓)	AP _p (↑)	mAP (↑)
Stable Diffusion [39]	✓	✓	✓					81.3	31.1	13.2
	✓							83.5	23.1	10.0
			✓					82.0	24.0	10.9
				✓				82.2	34.5	14.2
					✓			82.1	30.0	15.2
					✓	✓		82.0	33.5	16.7
					✓		✓	86.3	30.9	14.0
						✓	✓	84.1	23.5	10.5
						✓	✓	84.5	35.6	16.5
						✓		81.2	36.0	15.1
						✓	87.5	35.0	14.3	
							✓	85.0	25.6	9.8

Table A5. Ablation study of loss constraints on different layers.

Loss Scale	FID (↓)	AP _p (↑)	mAP (↑)
5	82.5	28.3	12.4
10	82.0	30.0	14.5
20	81.1	34.7	15.4
30	82.0	33.5	16.7
50	83.8	35.8	15.6
100	88.4	34.9	14.3
200	99.2	32.2	13.8
500	129.7	26.2	9.2

Table A6. Ablation study of the loss scale factor.

thus creating a prompt that consists only of padding tokens. Then, we use this prompt to generate images. Surprisingly, despite only generating from padding (*i.e.*, non-word) token embeddings, we observe that the generated images (Word Drop in Figure A3) closely follow both the semantics and the layout of the image generated from the full-text prompt. Thus, the figure clearly demonstrates that the padding tokens contain the information of the whole sentence. This further justifies why in forward guidance padding tokens cannot be ignored, *i.e.*, it would be insufficient to attempt to control selected word tokens only (main paper, Figure 5). In backward guidance, however, controlling the cross-attention maps of padding tokens is not necessary; this is now done by back-propagating and updating the latent, which subsequently changes the cross-attention maps of all tokens, even those that are not explicitly controlled.

Cross-Attention Maps of Special Tokens. During our experiments, we found that the cross-attention of the padding tokens has a strong connection to the foreground of the generated images. We illustrated this in Figure 4 (main paper), which shows that the cross-attention maps of padding tokens resemble saliency maps, while the cross-attention maps of the start tokens are mostly complementary to those of padding tokens (*i.e.*, they capture what can be considered as background). In Figure A4, we show more examples of the cross-attention maps of the *start* and

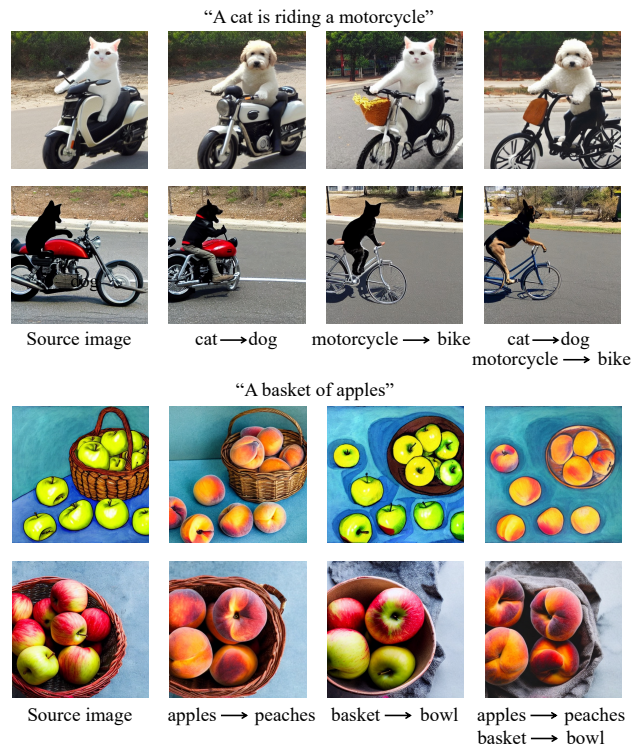


Figure A1. Each row has the same initial noise. We could see that even if we changed the object word in one sentence, the overall layout remains similar.

padding tokens. The captions are randomly taken from MSCOCO [28]. This figure further highlights the observation that cross-attention maps of these special tokens contain important semantic and spatial information. For example, in the first row, given “A short train traveling through a mountainous landscape” as the input prompt, the cross-attention map of the padding tokens aligns with the generated train and the start token focuses on the background of the generated image.

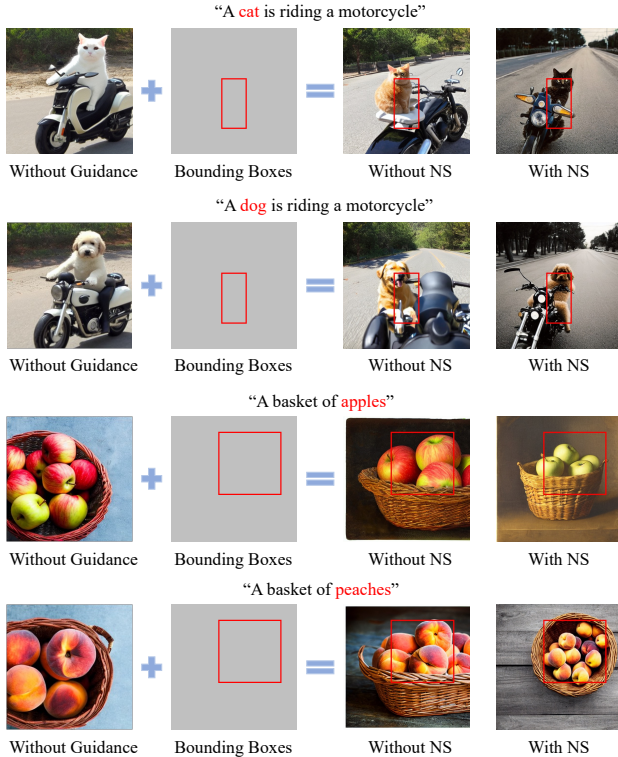


Figure A2. We qualitatively compare the generated results with and without noise selection (NS). The results show that with noise selection, our backward guidance achieves better layout control.

Base Model	NS	COCO 2014		Flickr30K		
		FID (↓)	mAP (↑)	FID (↓)	mAP (↑)	AP _P (↑)
Stable Diffusion	✗	74.4	33.6	82.0	33.5	16.7
Stable Diffusion	✓	73.3	35.7	78.9	35.6	17.9

Table A7. Ablation Study on Noise Selection (NS).

Layout Control with Only Padding Tokens. Motivated by the examples above, we perform backward guidance only on the cross-attention maps of padding tokens to control the spatial layout of all foreground objects simultaneously (as a group). Some examples are shown in Figure A5. This figure verifies our assumption that by guiding the cross-attention map of the padding tokens alone one can control the composition of the images at the foreground/background level.

F. More Examples.

More Examples under VISOR Protocol. We show more examples under the VISOR protocol in Figure A6 and Figure A7. Our method generates the correct spatial relationships as shown in the figures. There are also some failure cases, such as the last row in Figure A6. Our method

fails to generate both a fork and a carrot. This is an inherited problem from the Stable Diffusion model. However, in most cases, layout guidance helps generate *all* entities in the text prompt, even when the unguided Stable Diffusion fails (*e.g.*, as is often the case with atypical scene compositions), as well as conforming to a specific spatial arrangement.

More Image Editing Examples. We show more examples of real image editing in Figure A8. Specifically, we train for 500 steps to learn the embedding of $\langle * \rangle$ with text inversion and then 150 steps fine-tuning of the text encoder and denoiser network with Dreambooth. After finalizing the model, we perform inference with our backward guidance using different text prompts and user-specified bounding boxes. As shown in the figure, we manage to change the context, layout, and style of the given real image.

Text Prompt	Original Image	Word Drop	Text Prompt	Original Image	Word Drop
“A double decker bus driving down a street.”			“a close up of a hot dog on a table”		
“A cooked pizza on a silver plater with another in the background.”			“A fluffy black cat is laying on a bed.”		
“Several elephants walking together in a line near water.”			“a large giraffe is outside eating from a tree”		
“Sheep graze in a valley under a clear blue sky.”			“A stop sign with dirty edges at a cross walk of a street.”		
“A plate of food with bread, grape tomatoes, cheese, cucumbers and sauce on it.”			“A cup of coffee in a to-go cup and three pastries”		
“Closeup of various oranges and bananas in pile.”			“A hotel room with items strewn about it.”		
“a couple of bears that are leaning on a rock”			“A half eaten sandwich next to a partially eaten bowl of macaroni salad.”		

Figure A3. Generating images without “seeing” the full-text prompt. We replace the token embeddings for all words in each caption with their *padding* token embeddings (word drop). We observe that the generated images after word dropping exhibit similar semantics and layout to the images generated from the full-text prompt, suggesting that significant information about the image is contained in padding tokens.

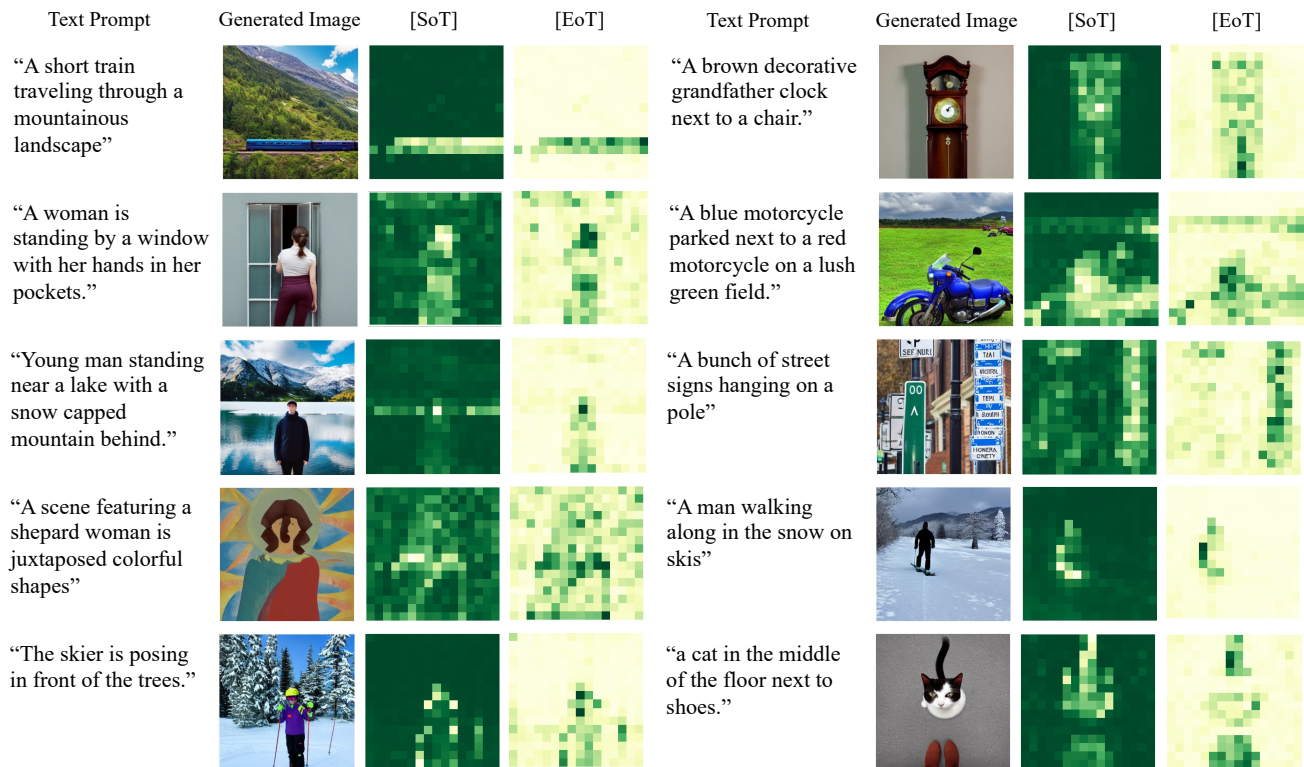


Figure A4. Visualization of cross-attention maps of start token ([SoT]) and padding tokens ([EoT]) at the final step of inference. Cross-attention maps are taken from the first cross-attention block of the up-sampling branch of U-Net and averaged over all attention heads.

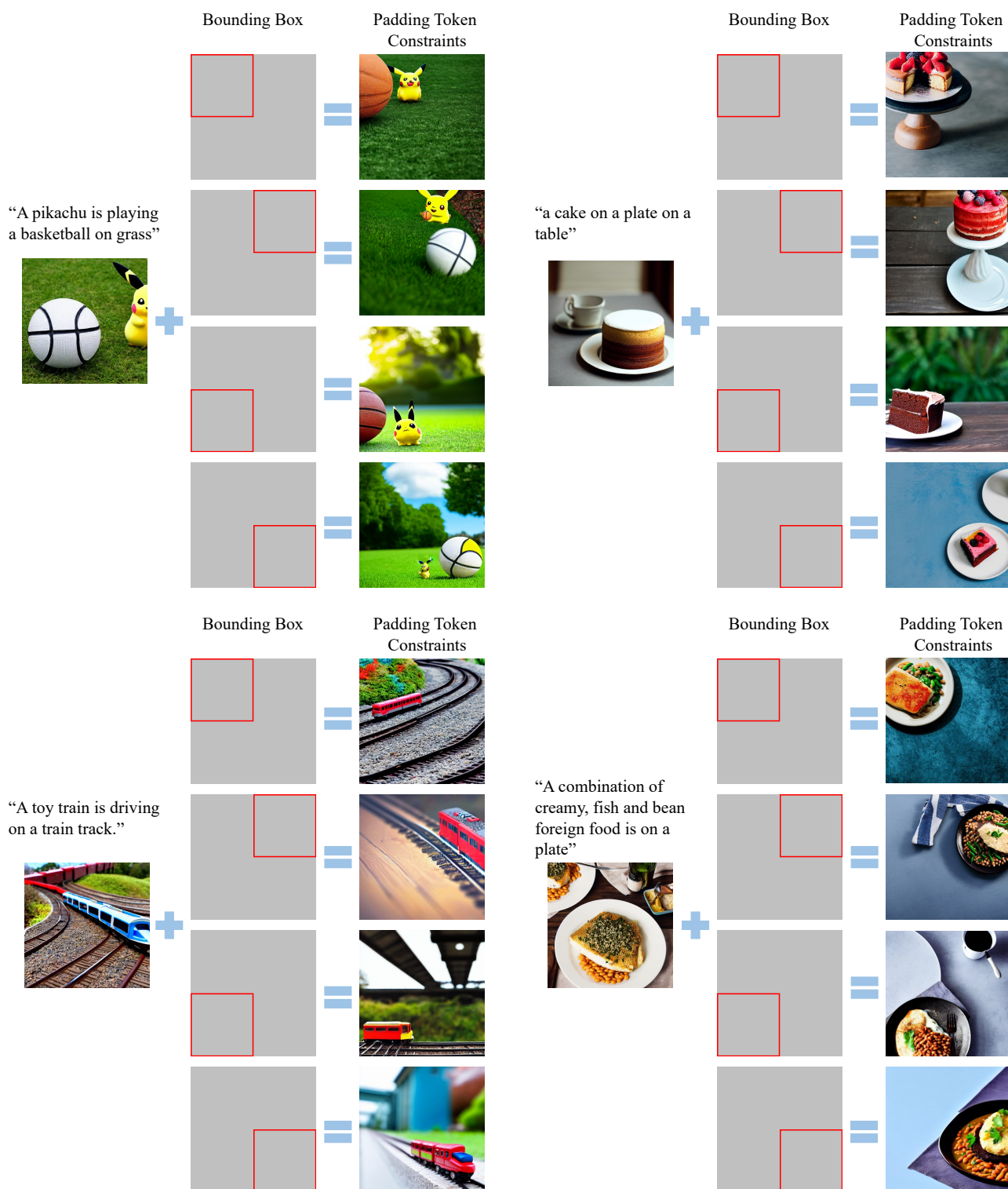


Figure A5. Backward guidance *only* on the *padding* tokens. We observe that the cross-attention of padding tokens typically represents the foreground of the generated image. Therefore, by spatially guiding the cross-attention maps that correspond to padding tokens, we can control the position of the foreground, which may include multiple objects (e.g., “pikachu” and “basketball”).

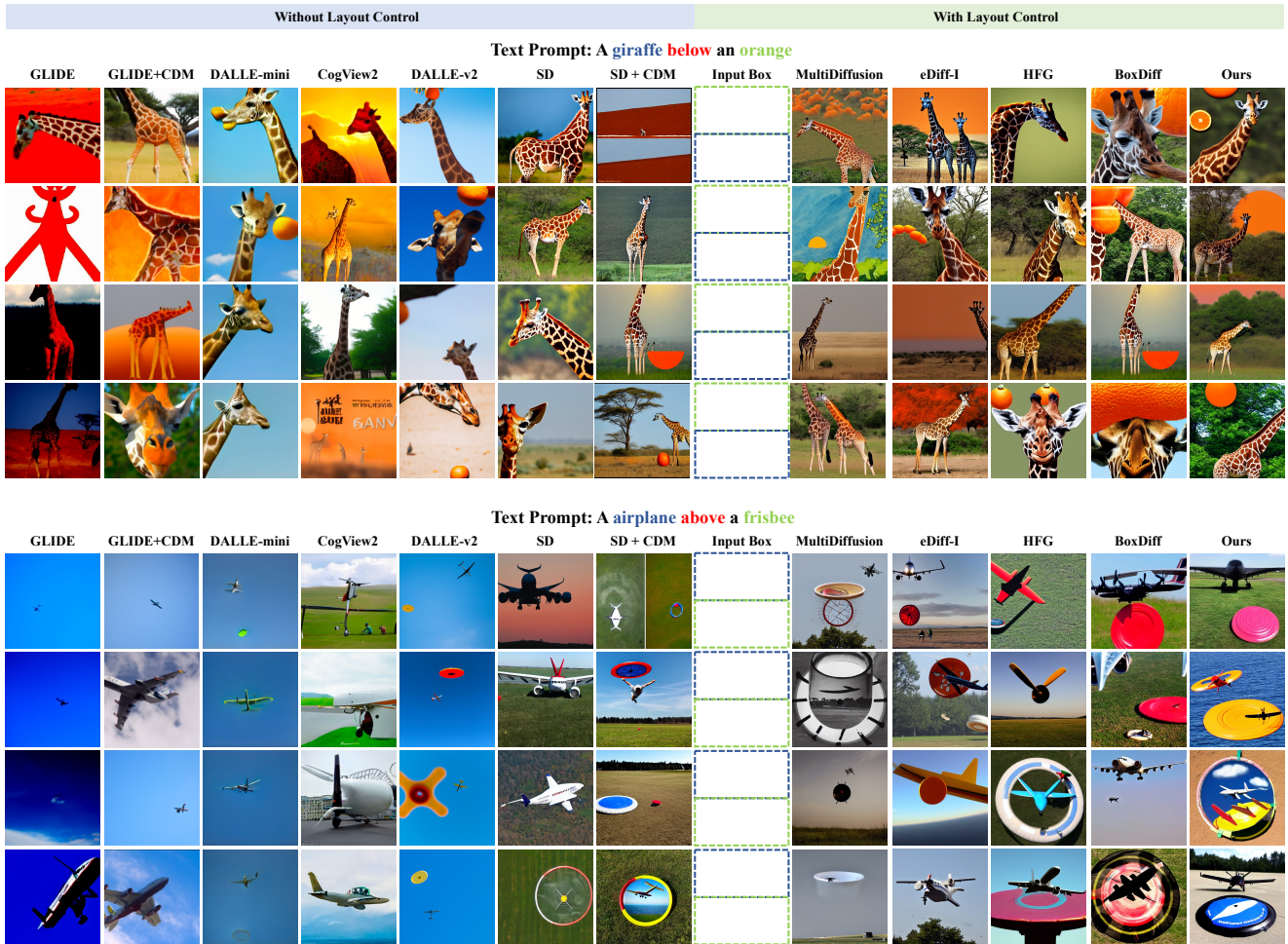


Figure A6. Qualitative comparison between different generative models. For each prompt, we generate four images. Some images of other models are from the demo website of [16].

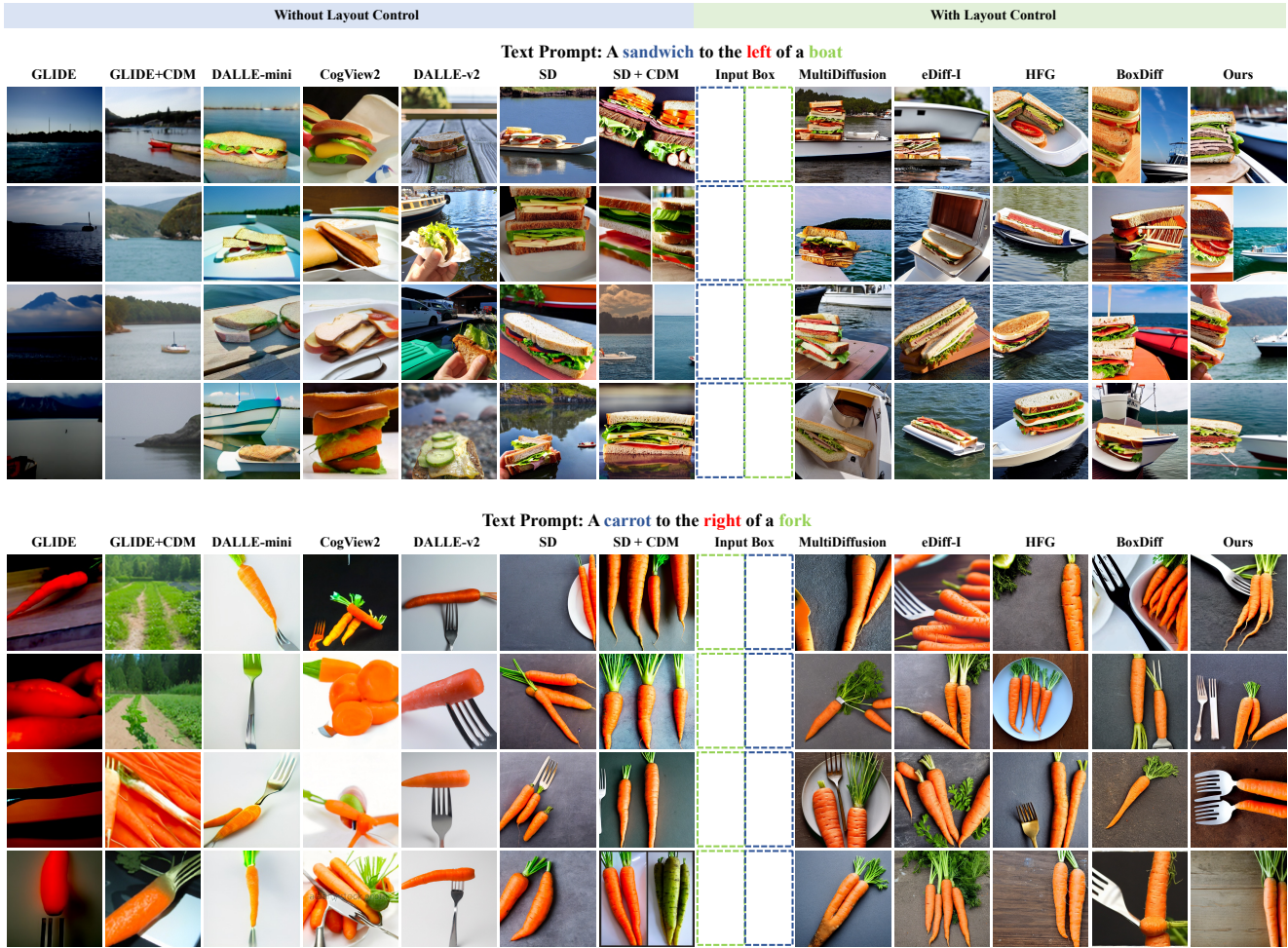


Figure A7. Qualitative comparison between different generative models. For each prompt, we generate four images. Some images of other models are from the demo website of [16].

Input Image



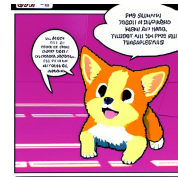
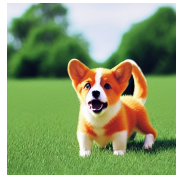
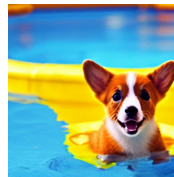
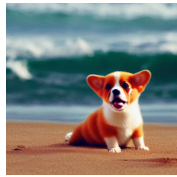
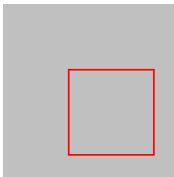
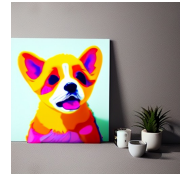
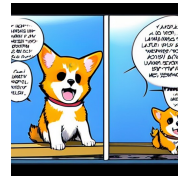
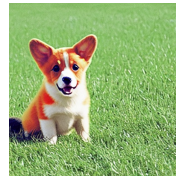
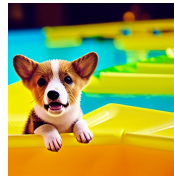
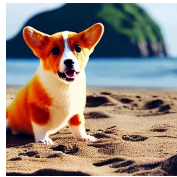
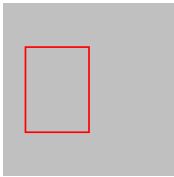
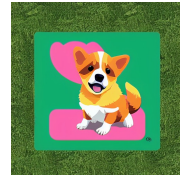
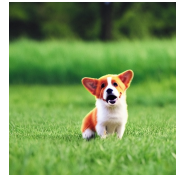
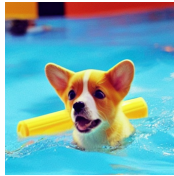
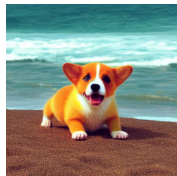
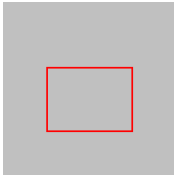
A $\langle * \rangle$ is standing on a beautiful beach.

A $\langle * \rangle$ is swimming in a pool.

A $\langle * \rangle$ is on grass.

A manga version of $\langle * \rangle$.

A painted version of $\langle * \rangle$.



Input Image



A $\langle * \rangle$ is standing on a beautiful beach.

A $\langle * \rangle$ is swimming in a pool.

A $\langle * \rangle$ is on grass.

A manga version of $\langle * \rangle$.

A painted version of $\langle * \rangle$.

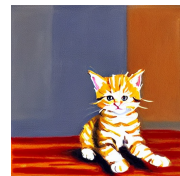
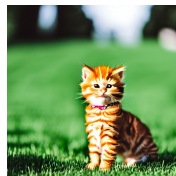
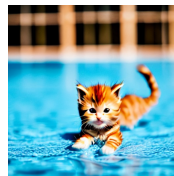
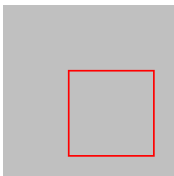
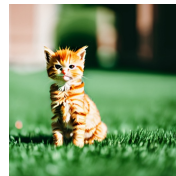
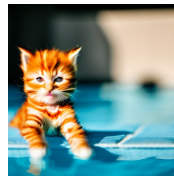
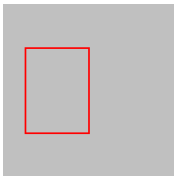
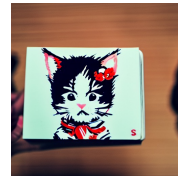
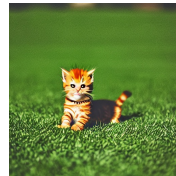
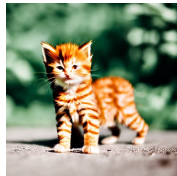
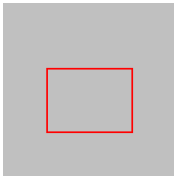


Figure A8. More examples of real image editing. $\langle * \rangle$ is the learned token that encodes the object in the real image.