

# GLIGEN: Open-Set Grounded Text-to-Image Generation

Yuheng Li<sup>1§</sup>, Haotian Liu<sup>1§</sup>, Qingyang Wu<sup>2</sup>, Fangzhou Mu<sup>1</sup>, Jianwei Yang<sup>3</sup>, Jianfeng Gao<sup>3</sup>,  
Chunyuan Li<sup>3¶</sup>, Yong Jae Lee<sup>1¶</sup>

<sup>1</sup>University of Wisconsin-Madison <sup>2</sup>Columbia University <sup>3</sup>Microsoft

<https://gligen.github.io/>

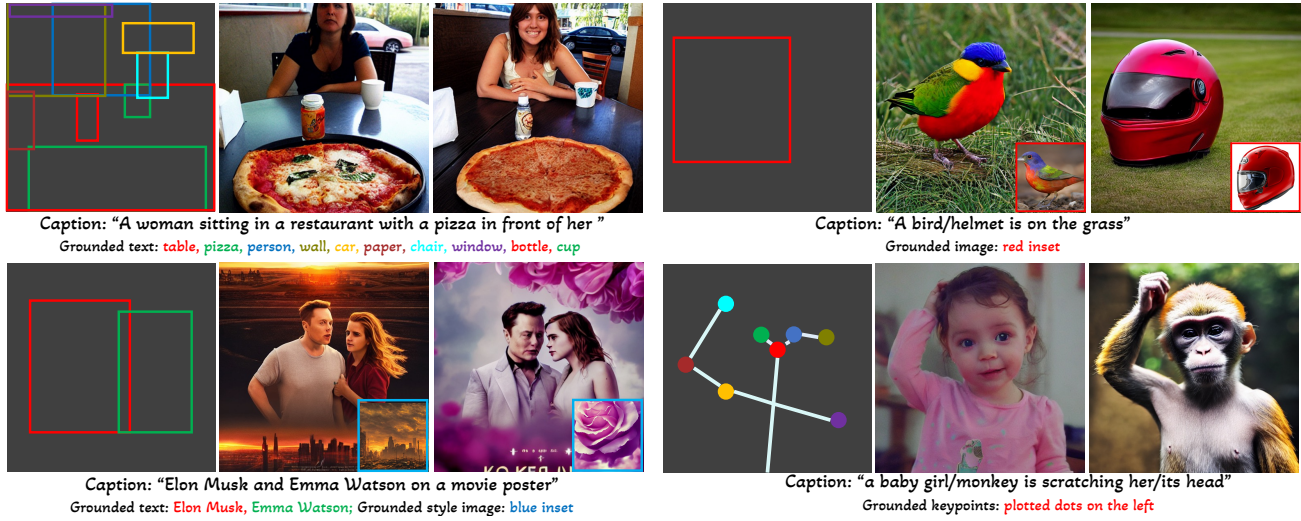


Figure 1. GLIGEN enables versatile grounding capabilities for a frozen text-to-image generation model, by feeding different grounding conditions. GLIGEN supports (a) text entity + box, (b) image entity + box, (c) image style and text + box, (d) text entity + keypoints. The generated examples for each scenario are shown in top-left, top-right, bottom-left, and bottom-right, respectively.

## Abstract

Large-scale text-to-image diffusion models have made amazing advances. However, the status quo is to use text input alone, which can impede controllability. In this work, we propose GLIGEN, **Grounded-Language-to-Image Generation**, a novel approach that builds upon and extends the functionality of existing pre-trained text-to-image diffusion models by enabling them to also be conditioned on grounding inputs. To preserve the vast concept knowledge of the pre-trained model, we freeze all of its weights and inject the grounding information into new trainable layers via a gated mechanism. Our model achieves open-world grounded text2img generation with caption and bounding box condition inputs, and the grounding ability generalizes well to novel spatial configurations and concepts. GLIGEN’s zero-shot performance on COCO and LVIS outperforms existing supervised layout-to-image baselines by a large margin.

§ Part of the work performed at Microsoft; ¶ Co-senior authors

## 1. Introduction

Image generation research has witnessed huge advances in recent years. Over the past couple of years, GANs [14] were the state-of-the-art, with their latent space and conditional inputs being well-studied for controllable manipulation [46, 58] and generation [27, 29, 45, 80]. Text conditional autoregressive [50, 72] and diffusion [49, 54] models have demonstrated astonishing image quality and concept coverage, due to their more stable learning objectives and large-scale training on web image-text paired data. These models have gained attention even among the general public due to their practical use cases (e.g., art design and creation).

Despite exciting progress, existing large-scale text-to-image generation models cannot be conditioned on other input modalities apart from text, and thus lack the ability to precisely localize concepts or use reference images to control the generation process. The current input, *i.e.*, natural language alone, restricts the way that information can

be expressed. For example, it is difficult to describe the precise location of an object using text, whereas bounding boxes / keypoints can easily achieve this, as shown in Figure 1. While conditional diffusion models [10, 51, 53] and GANs [26, 37, 46, 69] that take in input modalities other than text for inpainting, layout2img generation, *etc.*, do exist, they rarely combine those inputs for controllable text2img generation.

Moreover, prior generative models—regardless of the generative model family—are usually independently trained on each task-specific dataset. In contrast, in the recognition field, the long-standing paradigm has been to build a task-specific recognition model [32] by starting from a foundation model pretrained on large-scale image data [4, 16, 17] or image-text pairs [33, 48, 73]. Since diffusion models have been trained on billions of image-text pairs [51], a natural question is: *Can we build upon existing pretrained diffusion models and endow them with new conditional input modalities?* In this way, analogous to the recognition literature, we may be able to achieve better performance on other generation tasks due to the vast concept knowledge that the pretrained models have, while acquiring more controllability over existing text-to-image generation models.

With the above aims, we propose a method for providing new grounding conditional inputs to pretrained text-to-image diffusion models. As shown in Figure 1, we still retain the text caption as input, but also enable other input modalities such as bounding boxes for grounding concepts, grounding reference images, and grounding part keypoints. The key challenge is preserving the original vast concept knowledge in the pretrained model while learning to inject the new grounding information. To prevent knowledge forgetting, we propose to freeze the original model weights and add new trainable gated Transformer layers [65] that take in the new grounding input (*e.g.*, bounding box). During training, we gradually fuse the new grounding information into the pretrained model using a gated mechanism [1]. This design enables flexibility in the sampling process during generation for improved quality and controllability; for example, we show that using the full model (all layers) in the first half of the sampling steps and only using the original layers (without the gated Transformer layers) in the latter half can lead to generation results that accurately reflect the grounding conditions while also having high image quality.

In our experiments, we primarily study grounded text2img generation with bounding boxes, inspired by the recent scaling success of learning grounded language-image understanding models with boxes in GLIP [34]. To enable our model to ground open-world vocabulary concepts [32, 34, 74, 77], we use the same pre-trained text encoder (for encoding the caption) to encode each phrase associated with each grounded entity (*i.e.*, one phrase per bounding box) and feed the encoded tokens into the newly inserted

layers with their encoded location information. Due to the shared text space, we find that our model can generalize to unseen objects even when only trained on the COCO [41] dataset. Its generalization on LVIS [15] outperforms a strong fully-supervised baseline by a large margin. To further improve our model’s grounding ability, we unify the object detection and grounding data formats for training, following GLIP [34], as they provide complementary benefits: detection data is of larger quantity, while grounding data has a richer vocabulary. With larger training data, our model’s generalization is consistently improved.

**Contributions.** 1) We propose a new text2img generation method that endows new grounding controllability over existing text2img diffusion models. 2) By preserving the pre-trained weights and learning to gradually integrate the new localization layers, our model achieves open-world grounded text2img generation with bounding box inputs, *i.e.*, synthesis of novel localized concepts unobserved in training. 3) Our model’s zero-shot performance on layout2img tasks significantly outperforms the prior state-of-the-art, demonstrating the power of building upon large pretrained generative models for downstream tasks.

## 2. Related Work

**Large scale text-to-image generation models.** State-of-the-art models in this space are either autoregressive [13, 50, 67, 72] or diffusion [43, 49, 51, 54, 79]. Among autoregressive models, DALL-E [50] is one of the breakthrough works that demonstrates zero-shot abilities, while Parti [72] demonstrates the feasibility of scaling up autoregressive models. Diffusion models have also shown very promising results. DALL-E 2 [49] generates images from the CLIP [48] image space, while Imagen [54] finds the benefit of using pretrained language models. The concurrent Muse [6] demonstrates that masked modeling can achieve SoTA-level generation performance with higher inference speed. However, all of these models usually only take a caption as the input, which can be difficult for conveying other information such as the precise location of an object. Make-A-Scene [13] also incorporates semantic maps into its text-to-image generation, by training an encoder to tokenize semantic masks to condition the generation. However, it can only operate in a closed-set (of 158 categories), whereas our grounded entities can be open-world. A concurrent work eDiff-I [3] shows that by changing the attention map, one can generate objects that roughly follow a semantic map input. However, we believe our interface with boxes is simpler, and more importantly, our method allows other conditioning inputs such as keypoints, which are hard to manipulate through attention.

**Image generation from layouts.** Given bounding boxes labeled with object categories, the task is to generate a corresponding image [24, 39, 59–61, 70, 76], which is the reverse

task of object detection. Layout2Im [76] formulated the problem and combined a VAE object encoder, an LSTM [22] object fuser, and an image decoder to generate the image, using global and object-level adversarial losses [14] to enforce realism and layout correspondence. LostGAN [59, 60] generates a mask representation which is used to normalize features, taking inspiration from StyleGAN [28]. LAMA [39] improves the intermediate mask quality for better image quality. Transformer [64] based methods [24, 70] have also been explored. Critically, existing layout2image methods are closed-set, *i.e.*, they can only generate limited localized visual concepts observed in the training set such as the 80 categories in COCO. In contrast, our method represents the first work for *open-set* grounded image generation. A concurrent work ReCo [71] also demonstrates open-set abilities by building upon a pretrained Stable Diffusion model [51]. However, it finetunes the original model weights, which has the potential to lead to knowledge forgetting. Furthermore, it only demonstrates box grounding results whereas we also show image and keypoint grounding results.

**Other conditional image generation.** For GANs, various conditioning information have been explored; *e.g.*, text [63, 68, 78], box [59, 60, 76], semantic masks [36, 45], images [8, 38, 81]. For diffusion models, LDM [51] proposes a unified approach for conditional generation by injecting the condition via cross-attention layers. Palette [53] performs image-to-image tasks using diffusion models. These models are usually trained from scratch independently. In our work, we investigate how to build upon existing models pretrained on large-scale web data, to enable new open-set grounded image generation capabilities in a cost-effective manner.

### 3. Preliminaries on Latent Diffusion Models

Diffusion-based methods are one of the most effective model families for text2image tasks, among which latent diffusion model (LDM) [51] and its successor Stable Diffusion are the most powerful models publicly available to the research community. To reduce the computational costs of vanilla diffusion model training, LDM proceeds in two stages. The first stage learns a bidirectional mapping network to obtain the latent representation  $z$  of the image  $x$ . The second stage trains a diffusion model on the latent  $z$ . Since the first stage model produces a fixed bidirectional mapping between  $x$  and  $z$ , from hereon, we focus on the latent generation space of LDM for simplicity.

**Training Objective.** Starting from noise  $z_T$ , the model gradually produces less noisy samples  $z_{T-1}, z_{T-2}, \dots, z_0$ , conditioned on caption  $c$  at every time step  $t$ . To learn such a model  $f_\theta$  parameterized by  $\theta$ , for each step, the LDM training objective solves the denoising problem on latent

representations  $z$  of the image  $x$ :

$$\min_{\theta} \mathcal{L}_{\text{LDM}} = \mathbb{E}_{z, \epsilon \sim \mathcal{N}(0, \mathbf{I}), t} [\|\epsilon - f_\theta(z_t, t, c)\|_2^2], \quad (1)$$

where  $t$  is uniformly sampled from time steps  $\{1, \dots, T\}$ ,  $z_t$  is the step- $t$  noisy variant of input  $z$ , and  $f_\theta(*, t, c)$  is the  $(t, c)$ -conditioned denoising autoencoder.

**Network Architecture.** The core of the network architecture is how to encode the conditions, based on which a cleaner version of  $z$  is produced. (i) *Denoising Autoencoder.*  $f_\theta(*, t, c)$  is implemented via UNet [52]. It takes in a noisy latent  $z$ , as well as information from time step  $t$  and condition  $c$ . It consists of a series of ResNet [19] and Transformer [65] blocks. (ii) *Condition Encoding.* In the original LDM, a BERT-like [9] network is trained from scratch to encode each caption into a sequence of text embeddings,  $f_{\text{text}}(c)$ , which is fed into (1) to replace  $c$ . The caption feature is encoded via a fixed CLIP [48] text encoder in Stable Diffusion. Time  $t$  is first mapped to time embedding  $\phi(t)$ , then injected into the UNet. The caption feature is used in a cross attention layer within each Transformer block. The model learns to predict the noise, following (1).

With large-scale training, the model  $f_\theta(*, t, c)$  is well trained to denoise  $z$  based on the caption information only. Though impressive language-to-image generation results have been shown with LDM by pretraining on internet-scale data, it remains challenging to synthesize images where additional grounding input can be instructed, and is thus the focus of our paper.

## 4. Open-set Grounded Image Generation

### 4.1. Grounding Instruction Input

For grounded text-to-image generation, there are a variety of ways to ground an object via spatial conditioning. We denote as  $e$  the grounding entity described either through text or an example image, and as  $l$  the grounding spatial configuration described with *e.g.*, a bounding box or a set of keypoints. We define the instruction to a grounded text-to-image model as a composition of the caption and grounded entities:

$$\text{Instruction: } y = (c, e), \quad \text{with} \quad (2)$$

$$\text{Caption: } c = [c_1, \dots, c_L] \quad (3)$$

$$\text{Grounding: } e = [(e_1, l_1), \dots, (e_N, l_N)] \quad (4)$$

where  $L$  is the caption length, and  $N$  is the number of entities to ground. In this work, we primarily study using bounding box as the grounding spatial configuration  $l$ , because of its large availability and easy annotation for users. For the grounded entity  $e$ , we mainly focus on using text as its representation due to simplicity. We process both caption and grounding entities as input tokens to the diffusion model, as described in detail below.



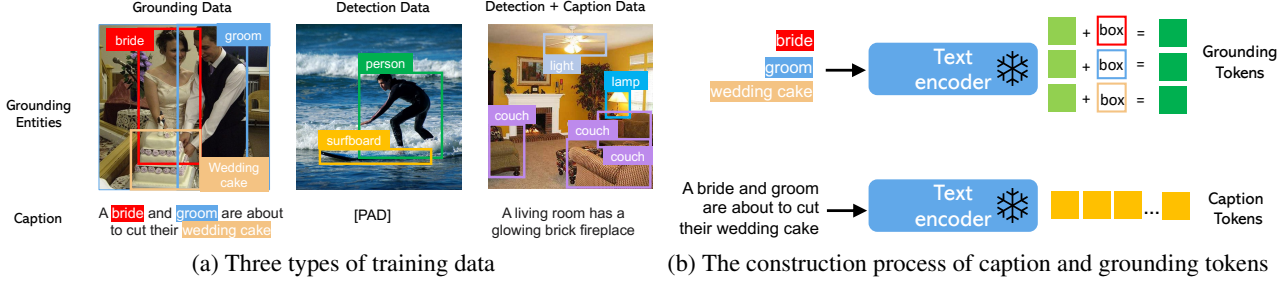


Figure 2. Illustration of training data and grounding instruction input. (a) For the grounding entities, we directly visualize the concept and bounding box information on the ground-truth images. The box is parameterized as normalized image coordinates, e.g. the person box is  $[0.37, 0.31, 0.92, 0.84]$ . (b) The first example in (a) is used to illustrate the token construction process.

**Caption Tokens.** The caption  $c$  is processed in the same way as in LDM. Specifically, we obtain the caption feature sequence (yellow tokens in Figure 2(b)) using  $\mathbf{h}^c = [h_1^c, \dots, h_L^c] = f_{\text{text}}(c)$ , where  $h_\ell^c$  is the contextualized text feature for the  $\ell$ -th word in the caption.

**Grounding Tokens.** For each grounded text entity denoted with a bounding box, we represent the location information as  $\mathbf{l} = [\alpha_{\min}, \beta_{\min}, \alpha_{\max}, \beta_{\max}]$  with its top-left and bottom-right coordinates. For the text entity  $e$ , we use the same pre-trained text encoder to obtain its text feature  $f_{\text{text}}(e)$  (light green token in Figure 2(b)), and then fuse it with its bounding box information to produce a grounding token (dark green token in Figure 2(b)):

$$\mathbf{h}^e = \text{MLP}(f_{\text{text}}(e), \text{Fourier}(\mathbf{l})) \quad (5)$$

where Fourier is the Fourier embedding [42], and  $\text{MLP}(\cdot, \cdot)$  is a multi-layer perceptron that first concatenates the two inputs across the feature dimension. The grounding token sequence is represented as  $\mathbf{h}^e = [h_1^e, \dots, h_N^e]$ .

**From Closed-set to Open-set.** Note that existing layout2img works only deal with a closed-set setting (e.g., COCO categories), as they typically learn a vector embedding  $\mathbf{u}$  per entity, to replace  $f_{\text{text}}(e)$  in (5). For a closed-set setting with  $K$  concepts, a dictionary of with  $K$  embeddings are learned,  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_K]$ . While this non-parametric representation works well in the closed-set setting, it has two drawbacks: (1) The conditioning is implemented as a dictionary look-up over  $\mathbf{U}$  in the evaluation stage, and thus the model can only ground the observed entities in the generated images, lacking the ability to generalize to ground new entities; (2) No word/phrase is ever utilized in the model condition, and the semantic structure [23] of the underlying language instruction is missing. In contrast, in our open-set design, since the noun entities are processed by the same text encoder that is used to encode the caption, we find that even when the localization information is limited to the concepts in the grounding training datasets, our model can still generalize to other concepts as we will show in our experiments.

**Training Data.** The training data for grounded image generation requires both text  $c$  and grounding entity  $e$  as the full

condition. In practice, we can relax the data requirement by considering a more flexible input, i.e. the three types of data shown in Figure 2(a). (i) *Grounding data.* Each image is associated with a caption describing the whole image; noun entities are extracted from the caption, and are labeled with bounding boxes. Since the noun entities are taken directly from the natural language caption, they can cover a much richer vocabulary which will be beneficial for open-world vocabulary grounded generation. (ii) *Detection data.* Noun entities are pre-defined closed-set categories (e.g., 80 object classes in COCO [41]). In this case, we choose to use a null caption token as introduced in classifier-free guidance [21] for the caption. The detection data is of larger quantity (millions) than the grounding data (thousands), and can therefore greatly increase overall training data. (iii) *Detection and caption data.* Noun entities are same as those in the detection data, and the image is described separately with a text caption. In this case, the noun entities may not exactly match those in the caption. For example, in Figure 2(a), the caption only gives a high-level description of the living room without mentioning the objects in the scene, whereas the detection annotation provides more fine-grained object-level details.

**Extensions to Other Grounding Conditions.** Note that the proposed grounding instruction in Eq (4) is in a general form, though our description thus far has focused on the case of using text as entity  $e$  and bounding box as  $\mathbf{l}$  (the major setting of this paper). To demonstrate the flexibility of the GLIGEN framework, we also study two additional representative cases which extend the use scenario of Eq (4).

- *Image Prompt.* While language allows users to describe a rich set of entities in an open-vocabulary manner, sometimes more abstract and fine-grained concepts can be better characterized by example images. To this end, one may describe entity  $e$  using an image, instead of language. We use an image encoder to obtain feature  $f_{\text{image}}(e)$  which is used in place of  $f_{\text{text}}(e)$  in Eq (5) when  $e$  is an image.
- *Keypoints.* As a simple parameterization method to specify the spatial configuration of an entity, bounding boxes ease the user-machine interaction interface by providing



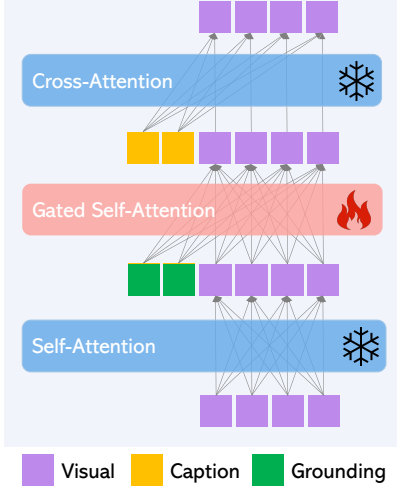


Figure 3. For a pretrained text2img model, the text features are fed into each cross-attention layer. A new gated self-attention layer is inserted to take in the new conditional localization information.

the height and width of the object layout only. One may consider richer spatial configurations such as keypoints for GLIGEN, by parameterizing  $\mathbf{l}$  in Eq (4) with a set of keypoint coordinates. Similar to encoding boxes, the Fourier embedding [42] can be applied to each keypoint location  $\mathbf{l} = [x, y]$ .

Figure 1 shows generated examples for these other grounding conditions. Please refer to the supp for more details.

## 4.2. Continual Learning for Grounded Generation

Our goal is to endow new spatial grounding capabilities to existing large language-to-image generation models. Large diffusion models have been pre-trained on web-scale image-text to gain the required knowledge for synthesizing realistic images based on diverse and complex language instructions. Due to the high pre-training cost and excellent performance, it is important to retain such knowledge in the model weights while expanding the new capability. Hence, we consider to lock the original model weights, and gradually adapt the model by tuning new modules.

**Gated Self-Attention.** We denote  $\mathbf{v} = [v_1, \dots, v_M]$  as the visual feature tokens of an image. The original Transformer block of LDM consists of two attention layers: The self-attention over the visual tokens, followed by cross-attention from caption tokens. By considering the residual connection, the two layers can be written:

$$\mathbf{v} = \mathbf{v} + \text{SelfAttn}(\mathbf{v}) \quad (6)$$

$$\mathbf{v} = \mathbf{v} + \text{CrossAttn}(\mathbf{v}, \mathbf{h}^c) \quad (7)$$

We freeze these two attention layers and add a new gated self-attention layer to enable the spatial grounding ability;

see Figure 3. Specifically, the attention is performed over the concatenation of visual and grounding tokens  $[\mathbf{v}, \mathbf{h}^g]$ :

$$\mathbf{v} = \mathbf{v} + \beta \cdot \tanh(\gamma) \cdot \text{TS}(\text{SelfAttn}([\mathbf{v}, \mathbf{h}^g])) \quad (8)$$

where  $\text{TS}(\cdot)$  is a token selection operation that considers visual tokens only, and  $\gamma$  is a learnable scalar which is initialized as 0.  $\beta$  is set as 1 during the entire training process and is only varied for scheduled sampling during inference (introduced below) for improved quality and controllability. Note that (8) is injected in between (6) and (7). Intuitively, the gated self-attention in (8) allows visual features to leverage bounding box information, and the resulting grounded features are treated as a residual, whose gate is initially set to 0 (due to  $\gamma$  being initialized as 0). This also enables more stable training. Note that a similar idea is used in Flamingo [1]; however, it uses gated cross-attention, which leads to worse performance in our case, possibly due to the lack of position embeddings for the visual features in the pretrained diffusion model.

**Learning Procedure.** We adapt the pre-trained model such that grounding information can be injected while all the original components remain intact. By denoting the new parameters in all gated self-attention layers as  $\theta'$ , we use the original denoising objective as in (1) for model continual learning, based on the grounding instruction input  $\mathbf{y}$ :

$$\min_{\theta'} \mathcal{L}_{\text{Grounding}} = \mathbb{E}_{\mathbf{z}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t} [\|\epsilon - f_{\{\theta, \theta'\}}(\mathbf{z}_t, t, \mathbf{y})\|_2^2]. \quad (9)$$

Why should the model try to use the new grounding information? Intuitively, predicting the noise that was added to a training image in the reverse diffusion process would be easier if the model could leverage the external knowledge about each object’s location. Thus, in this way, the model learns to use the additional localization information while retaining the pre-trained concept knowledge.

**A Versatile User Interface.** Once the model is well trained, our design of disentangling the caption and grounding inputs supports a versatile interface. Not only do we allow a user to ground entities that exist in the caption input, but objects can also be freely added in the desired locations without being mentioned in the caption input (see the pizza example in Figure 1). For a pure text-based diffusion model, a user would have to clumsily describe all the objects in the caption, while also specifying their precise locations, which can be difficult to do with language alone.

**Scheduled Sampling in Inference.** The standard inference scheme of GLIGEN is to set  $\beta = 1$  in (8), and the entire diffusion process is influenced by the grounding tokens. This constant  $\beta$  sampling scheme provides overall

good performance in terms of both generation and grounding, but sometimes generates lower quality images compared with the original text2img models (e.g., as Stable Diffusion is finetuned on high aesthetic scored images). To strike a better trade-off between generation and grounding for GLIGEN, we propose a scheduled sampling scheme. As we freeze the original model weights and add new layers to inject new grounding information in training, there is flexibility during inference to schedule the diffusion process to either use both the grounding and language tokens or use only the language tokens of the original model at anytime, by setting different  $\beta$  values in (8). Specifically, we consider a two-stage inference procedure, divided by  $\tau \in [0, 1]$ . For a diffusion process with  $T$  steps in total, one can set  $\beta$  to 1 at the beginning  $\tau * T$  steps, and set  $\beta$  to 0 for the remaining  $(1 - \tau) * T$  steps:

$$\beta = \begin{cases} 1, & t \leq \tau * T \quad \# \text{ Grounded inference stage} \\ 0, & t > \tau * T \quad \# \text{ Standard inference stage} \end{cases} \quad (10)$$

The major benefit of scheduled sampling is improved visual quality as the rough concept location and outline are decided in the early stages, followed by fine-grained details in later stages. It also allows us to extend the model trained in one domain (human keypoint) to other domains (monkey, cartoon characters) as shown in Figure 1.

## 5. Experiments

We evaluate our model’s grounded text2img generation in both the closed-set and open-set settings, ablate its components, and show extensions to image prompt and keypoint grounded generation. We conduct our main quantitative experiments by building upon a pretrained LDM on LAION [55], unless stated otherwise.

### 5.1. Closed-set Grounded Text2Img Generation

We first evaluate the generation quality and grounding accuracy of our model in a closed-set setting. For this, we train and evaluate on the COCO2014 [41] dataset, which is a standard benchmark used in the text2img literature [49, 54, 63, 68, 80], and evaluate how the different types of grounding instructions impact our model’s performance.

**Grounding instructions.** We use the following grounding instructions to train our model: 1) COCO2014D: Detection Data. There are no caption annotations so we use a null caption input [21]. Detection annotations are used as noun-entities. 2) COCO2014CD: Detection + Caption Data. Both caption and detection annotations are used. Note that the noun entities may not always exist in the caption. 3) COCO2014G: Grounding Data. Given the caption annotations, we use GLIP [34], which detects the caption’s noun entities in the image, to get pseudo box labels.

Model	Generation: FID ( $\downarrow$ )		Grounding: YOLO ( $\uparrow$ ) AP/AP <sub>50</sub> /AP <sub>75</sub>
	Fine-tuned	Zero-shot	
CogView [11]	-	27.10	-
KNN-Diffusion [2]	-	16.66	-
DALL-E 2 [49]	-	10.39	-
Imagen [54]	-	7.27	-
Re-Imagen [7]	5.25	6.88	-
Parti [72]	3.20	7.23	-
LAFITE [80]	8.12	26.94	-
LAFITE2 [78]	4.28	8.42	-
Make-a-Scene [13]	7.55	11.84	-
NÜWA [67]	12.90	-	-
Frido [12]	11.24	-	-
XMC-GAN [75]	9.33	-	-
AttnGAN [68]	35.49	-	-
DF-GAN [63]	21.42	-	-
Obj-GAN [35]	20.75	-	-
LDM [51]	-	12.63	-
LDM*	5.91	11.73	0.6 / 2.0 / 0.3
GLIGEN (COCO2014CD)	5.82	-	21.7 / 39.0 / 21.7
GLIGEN (COCO2014D)	5.61	-	<b>24.0 / 42.2 / 24.1</b>
GLIGEN (COCO2014G)	6.38	-	11.2 / 21.2 / 10.7

Table 1. Evaluation of image quality and correspondence to layout on COCO2014 val-set. All numbers are taken from corresponding papers, LDM\* is our COCO fine-tuned LDM baseline. Here GLIGEN is built upon LDM.

**Baselines.** Baseline models are listed in Table 1. Among them, we also finetune an LDM [51] pretrained on LAION 400M [55] on COCO2014 with its caption annotations, which we denote as LDM\*. The text2img baselines, as they cannot be conditioned on box inputs, are trained on COCO2014C: Caption Data.

**Evaluation metrics.** We use the captions and/or box annotations from 30K randomly sampled images to generate 30K images for evaluation. We use *FID* [20] to evaluate image quality. To evaluate grounding accuracy (*i.e.* correspondence between the input bounding box and generated entity), we use the *YOLO score* [40]. Specifically, we use a pretrained YOLO-v4 [5] to detect bounding boxes on the generated images and compare them with the ground truth boxes using average precision (AP). Since prior text2img methods do not support taking box annotations as input, it is not fair to compare with them on this metric. Thus, we only report numbers for the fine-tuned LDM as a reference.

**Results.** Table 1 shows the results. First, we see that the image synthesis quality of our approach, as measured by FID, is better than most of the state-of-the-art baselines due to rich visual knowledge learned in the pretraining stage. Next, we find that all three grounding instructions lead to comparable FID to that of the LDM\* baseline, which is finetuned on COCO2014 with caption annotations. Our model trained using detection annotation instructions (COCO2014D) has the overall best performance. However, when we evaluate this model on COCO2014CD instructions, we find that it has worse performance (FID: 8.2) – its ability to understand real captions may be limited as it is only trained with the null caption. For the model trained with GLIP grounding instructions (COCO2014G), we actually evaluate it using the COCO2014CD instructions since we need to compute

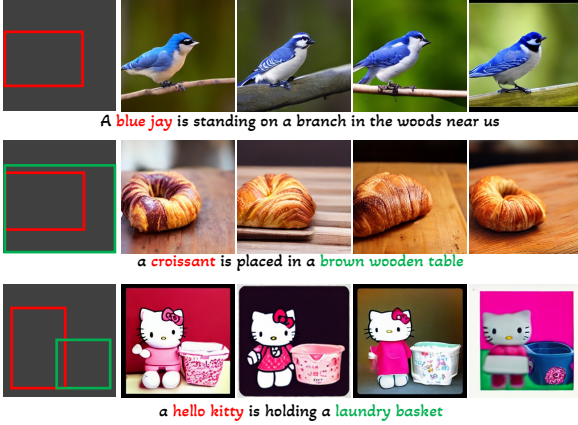


Figure 4. Our model can generalize to open-world concepts even when only trained using localization annotation from COCO.

the YOLO score which requires ground-truth detection annotations. Its slightly worse FID may be attributed to its learning from GLIP pseudo-labels. The same reason can explain its low YOLO score (*i.e.*, the model did not see any ground-truth detection annotations during training).

Overall, this experiment shows that: 1) Our model can successfully take in boxes as an additional condition while maintaining image generation quality. 2) All grounding instruction types are useful, which suggests that combining their data together can lead to complementary benefits.

**Ablation on gated self-attention.** Our approach uses gated self-attention to absorb the grounding instruction. We can also consider gated cross-attention [1], where the query is the visual feature, and the keys and values are produced using the grounding condition. We ablate this design on COCO2014CD data, and find that it leads to similar FID: 5.8, but worse YOLO AP: 16.6 (compared to 21.7 for self-attention in Table 1). This shows the necessity of information sharing among the visual tokens, which exists in self-attention but not in cross-attention.

**Ablation on null caption.** We choose to use the null caption when we only have detection annotations (COCO2014D). An alternative scheme is to simply combine all noun entities into a sentence; *e.g.*, if there are two cats and a dog in an image, then the pseudo caption can be: “cat, cat, dog”. In this case, the FID becomes worse and increases to 7.40 from 5.61 (null caption). This is likely due to the pretrained text encoder never having encountered this type of unnatural caption during LDM training. A solution would be to finetune the text encoder or design a better prompt, but this is not the focus of our work.

**Comparison to Layout2Img generation methods.** Thus far, we have seen that our model correctly learns to use the grounding condition. But how accurate is it compared to methods that are specifically designed for layout2img generation? To answer this, we train our model on COCO2017D,

Model	FID ( $\downarrow$ )	YOLO score (AP/AP <sub>50</sub> /AP <sub>75</sub> ) ( $\uparrow$ )
LostGAN-V2 [60]	42.55	9.1 / 15.3 / 9.8
OCGAN [62]	41.65	-
HCSS [25]	33.68	-
LAMA [40]	31.12	13.40 / 19.70 / 14.90
TwFA [69]	22.15	- / 28.20 / 20.12
GLIGEN-LDM	<b>21.04</b>	<b>22.4 / 36.5 / 24.1</b>

Table 2. Image quality and correspondence to layout are compared with baselines on COCO2017 val-set.

Model	Training data	AP	AP <sub>r</sub>	AP <sub>c</sub>	AP <sub>f</sub>
LAMA [40]	LVIS	2.0	0.9	1.3	3.2
GLIGEN-LDM	COCO2014CD	6.4	5.8	5.8	7.4
GLIGEN-LDM	COCO2014D	4.4	2.3	3.3	6.5
GLIGEN-LDM	COCO2014G	6.0	4.4	6.1	6.6
GLIGEN-LDM	GoldG,O365	10.6	5.8	9.6	13.8
GLIGEN-LDM	GoldG,O365,SBU,CC3M	11.1	9.0	9.8	13.4
GLIGEN-Stable	GoldG,O365,SBU,CC3M	10.8	8.8	9.9	12.6
Upper-bound	-	25.2	19.0	22.2	31.2

Table 3. GLIP-score on LVIS validation set. Upper-bound is provided by running GLIP on real images scaled to  $256 \times 256$ .

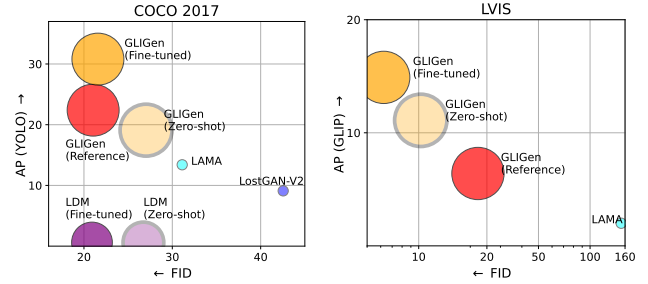


Figure 5. Performance comparison measured by image generation and grounding quality on COCO2017 (left) and LVIS (right) datasets. GLIGEN is built upon LDM, and continually pre-trained on the joint data of GoldG, O365, SBU, and CC3M. GLIGEN (Reference) is pre-trained on COCO/LVIS only. The circle size indicates the model size.

which only has detection annotations. We use the 2017 splits (instead of 2014 as before), as it is the standard benchmark in the layout2img literature. In this experiment, we use the exact same annotation as all layout2img baselines.

Table 2 shows that we achieve the state-of-the-art performance for both image quality and grounding accuracy. We believe the core reason is because previous methods train their model from scratch, whereas we build upon a large-scale pretrained generative model with rich visual semantics. Qualitative comparisons are in the supp. We also scale up our training data (discussed later) and pretrain a model on this dataset. Figure 5 left shows this model’s zero-shot and finetuned results.

## 5.2. Open-set Grounded Text2Img Generation

**COCO-training model.** We first take GLIGEN trained only with the grounding annotations of COCO (COCO2014CD), and evaluate whether it can generate grounded entities beyond the COCO categories. Figure 4 shows qualitative results, where GLIGEN can ground new concepts such as “blue jay”, “croissant” or ground object attributes





Figure 6. **Inpainting results.** Existing text2img diffusion models may generate objects that do not tightly fit the masked box or miss an object if the same object already exists in the image.

	1%-3%	5%-10%	30%-50%
LDM [51]	25.9	23.4	14.6
GLIGEN-LDM	<b>29.7</b>	<b>30.9</b>	<b>25.6</b>
Upper-bound	41.7	43.4	45.0

Table 4. Inpainting results (YOLO AP) for different size of objects.

such as “brown wooden table”, beyond the training categories. We hypothesize this is because the gated self-attention of GLIGEN learns to re-position the visual features corresponding to the grounding entities in the caption for the ensuing cross-attention layer, and gains generalization ability due to the shared text spaces in these two layers.

We also quantitatively evaluate our model’s zero-shot generation performance on LVIS [15], which contains 1203 long-tail object categories. We use GLIP to predict bounding boxes from the generated images and calculate AP, thus we name it as *GLIP score*. We compare to a state-of-the-art model designed for the layout2img task: LAMA [40]. We train LAMA using the official code on the LVIS training set (in a fully-supervised setting), whereas we directly evaluate our model in a *zero-shot task transfer* manner, by running inference on the LVIS val set without seeing any LVIS labels. Table 3 (first 4 rows) shows the results. Surprisingly, even though our model is only trained on COCO annotations, it outperforms the supervised baseline by a large margin. This is because the baseline, which is trained from scratch, struggles to learn from limited annotations (many of the rare classes in LVIS have fewer than five training samples). In contrast, our model can take advantage of the pretrained model’s vast concept knowledge.

**Scaling up the training data.** We next study our model’s open-set capability with much larger training data. Specifically, we follow GLIP [34] and train on Object365 [56] and GoldG [34], which combines two grounding datasets: Flickr [47] and VG [31]. We also use CC3M [57] and SBU [44] with grounding pseudo-labels generated by GLIP.

Table 3 shows the data scaling results. As we scale up the training data, our model’s zero-shot performance increases, especially for rare concepts. We also try to finetune the model pretrained on our largest dataset on LVIS and demon-

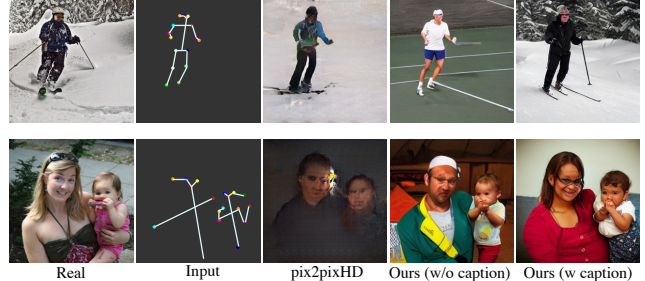


Figure 7. **Keypoint results.** Our model generates higher quality images conditioned on keypoints, and it allows to use caption to specify details such as scene or gender.

Model	FID	AP	AP <sub>50</sub>	AP <sub>75</sub>
pix2pixHD [66]	142.4	15.8	33.7	13.0
GLIGEN (w/o caption)	31.02	<b>31.8</b>	<b>53.5</b>	<b>31.0</b>
GLIGEN (w caption)	<b>27.34</b>	31.5	52.9	<b>31.0</b>
Upper-bound	-	62.4	75.0	65.9

Table 5. Conditioning with Human Keypoints evaluated on COCO2017 validation set. Upper-bound is calculated on real images scaled to  $256 \times 256$ .

strate its performance on Figure 5 right. To demonstrate the generality of our method, we also train our model based on the Stable Diffusion model checkpoint using the largest data. We show some qualitative examples in Figure 8 using this model. Our model gains the grounding ability compared to vanilla Stable Diffusion. We also notice that Stable Diffusion model may overlook certain objects (“umbrella” in the second example) due to its use of the CLIP text encoder which tends to focus on global scene properties, and may ignore object-level details [3]. It also struggles to generate spatially counterfactual concepts. By explicitly injecting entity information through grounding tokens, our model can improve the grounding ability in two ways: the referred objects are more likely to appear in the generated images, and the objects reside in the specified spatial location.

### 5.3. Inpainting Comparison

Like other diffusion models, GLIGEN can also work for the inpainting task by replacing the known region with a sample from  $q(z_t|z_0)$  after each sampling step, where  $z_0$  is the latent representation of an image [51]. One can ground text descriptions to missing regions, as shown in Figure 6. In this setting, however, one may wonder, can we simply use a vanilla text-to-image diffusion model such Stable Diffusion or DALL E2 to fill the missing region by providing the object name as the caption? What are the benefits of having extra grounding inputs in such cases? To answer this, we conduct the following experiment on the COCO dataset: for each image, we randomly mask one object. We then let the model inpaint the missing region. We choose the missing object with three different size ratios with respect to the image: small (1%-3%), median (5%-10%), and large (30%-50%). 5000 images are used for each case.

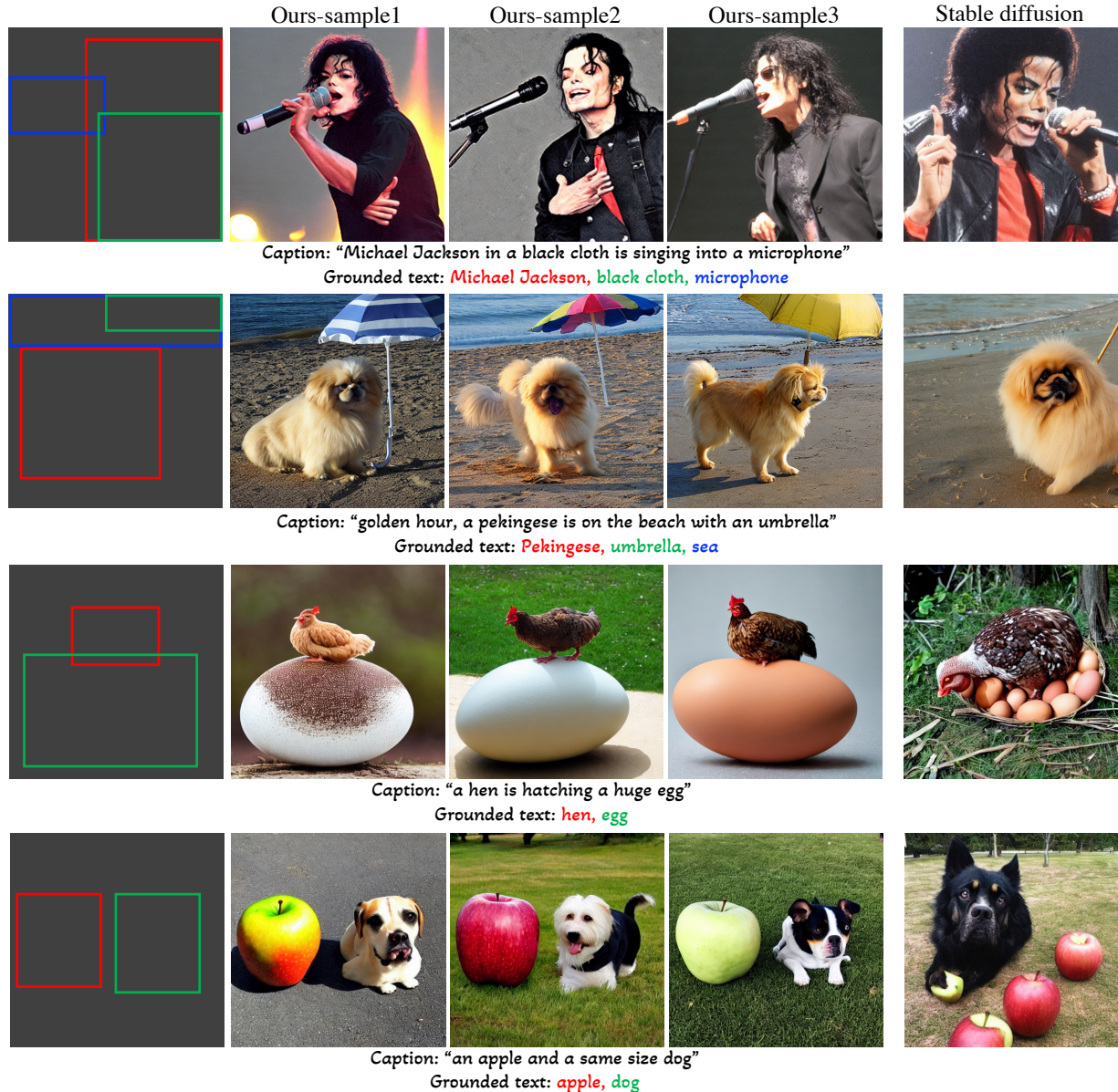


Figure 8. Grounded text2image generation. The baseline lacks grounding ability and can also miss objects e.g. “umbrella” in a sentence with multiple objects due to CLIP text space, and it also struggles to generate spatially counterfactual concepts.

Table 4 demonstrates that our inpainted objects more tightly occupy the missing region (box) compared to the baselines. Fig. 6 provides examples to visually compare the inpainting results (we use Stable Diffusion for better quality). The first row shows that baselines’ generated objects do not follow the provided box. The second row shows that when the missing category is already present in the image, they may ignore the caption. This is understandable as baselines are trained to generate a *whole* image following the caption. Our method may be more favorable for editing applications, where a user might want to generate an object that fully fits the missing region or add an instance of a class that already exists in the image.

#### 5.4. Keypoints Grounding

Although we have thus far demonstrated results with bounding boxes, our approach has flexibility in the grounding condition that it can use for generation. To demonstrate this, we next evaluate our model with another type of grounding condition: human keypoints. We use the COCO2017 dataset; details of the tokenization process for keypoints can be found in the supp. We compare with pix2pixHD [66], a classic image-to-image translation model. Since pix2pixHD does not take captions as input, we train two variants of our model: one uses COCO captions, the other does not. In the latter case, null caption is used as input to the cross-attention layer for a fair comparison.





Figure 9. **Image grounded generation** (top) where images can provide more fine-grained details than the text in the caption. **Text and image grounded generation** (bottom). Text is grounded using the red bounding box, and image is used as style reference.

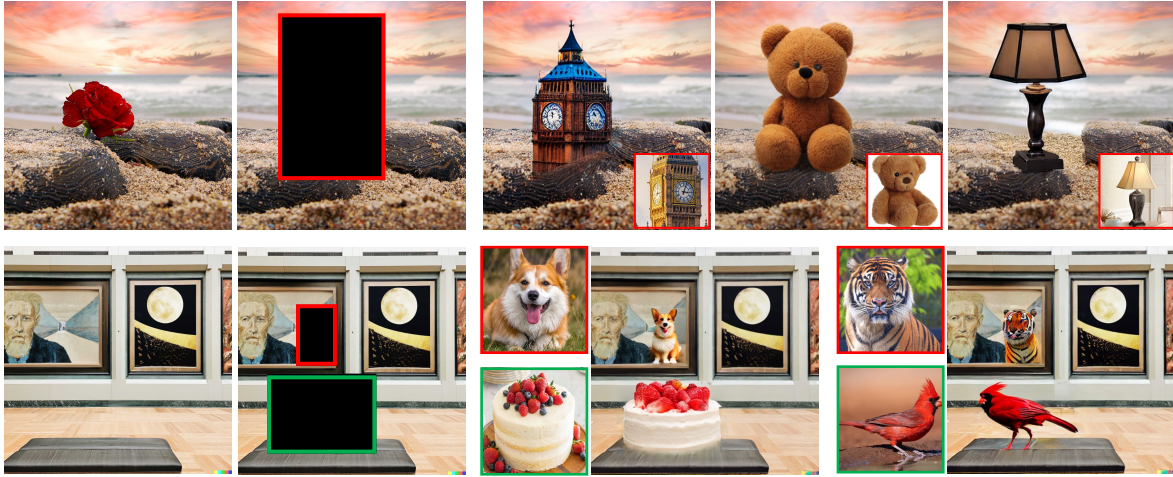


Figure 10. **Image grounded Inpainting**. One can use reference images to ground holes they want to fill in.

Fig. 7 shows the qualitative comparison. Clearly, our method generates much better image quality. For our model trained with captions, we can also specify other details such as the scene ("A person is skiing down a snowy hill") or person's gender ("A woman is holding a baby"). These two inputs complement each other and can enrich a user's controllability for image creation. We measure keypoint correspondence (similar to the YOLO score for boxes) by running a MaskRCNN [18] keypoint detector on the generated images. Both of our model variants produce similar results; see Table 5.

## 5.5. Image Grounding

**Image grounded generation.** One can also use a reference image to represent a grounded entity as discussed pre-

viously. Fig. 9 top row shows qualitative results, which demonstrate that the visual feature can complement details that are hard to describe by language, such as the style and shape of cars.

**Text and image grounded generation.** Besides using either text or image to represent a grounded entity, one can also keep both representations in one model for more creative generation. Figure 9 second row shows text grounded generation with style / tone transfer. Here we ground the text ("a brick house") with the red bounding box. For the style reference image, we find that grounding it to an image corner (green bounding box) or its edge is sufficient. Since the model needs to generate a harmonious style for the entire image, we hypothesize the self-attention layers may broadcast this information to all pixels, thus leading to consistent



style for the entire image.

**Image grounded inpainting.** As we previously demonstrated, one can ground text to missing region for inpainting, one can also ground reference images to missing regions. Figure 10 shows inpainting results grounded on reference images. To remove boundary artifacts, we follow GLIDE [43], and modify the first conv layer by adding 5 extra channels (4 for  $z_0$  and 1 for inpainting mask) and make them trainable with the new added layers.

## 5.6. Scheduled Sampling

As stated in Eq. (8) and Eq. (10), we can schedule inference time sampling by setting  $\beta$  to 1 (use extra grounding information) or 0 (reduce to the original pretrained diffusion model). This can make our model exploit different knowledge at different stages.

Fig. 11 qualitatively shows the benefits of our scheduled sampling for our model built upon Stable Diffusion. The images in the same row share the same noise and conditional input. The first row shows that scheduled sampling can be used to improve image quality, as the original Stable Diffusion model is trained with high quality images. The first 20% steps usually are sufficient for setting the overall structure of large objects, and the original Stable Diffusion model can then complete the remaining sampling process with its high-quality prior. The second row shows a generation example by our model trained with COCO human keypoint annotations. Since this model is purely trained with human keypoints, the final result is biased towards generating a human even if a different object (i.e., robot) is specified in the caption. However, by using scheduled sampling, we can extend this model to generate other objects with a human-like shape.

We do notice that scheduled sampling can decrease the correspondence to box or keypoint if we set the  $\tau$  to be too small. To quantitatively measure this, we evaluate the GLIP score of the generated images on the LVIS dataset (similar to Table 3) and we use GLIGEN based on the Stable Diffusion model. The GLIP AP for  $\tau = 0.0$  (vanilla Stable Diffusion),  $\tau = 0.2$ ,  $\tau = 0.3$ ,  $\tau = 0.5$ ,  $\tau = 0.8$  and  $\tau = 1$  (GLIGEN) are: 0.3, 2.4, 4.7, 8.9, 10.8. Note that to present the GLIGEN own performance, we set  $\tau = 1$ , i.e., without scheduled sampling, for all results reported in the previous sections.

## 6. Conclusion

We proposed GLIGEN for expanding pretrained text2img diffusion models with grounding ability, and demonstrated open-world generalization using bounding boxes as the grounding condition. Our method is simple and effective, and can be easily extended to other conditions *e.g.*, keypoints and reference images. One limitation we noticed is that the generated style or aesthetic distribution can shift after adding the new gated self-attention layers (*e.g.*, the model sometimes struggles to generate graphics style images when  $\tau$  is

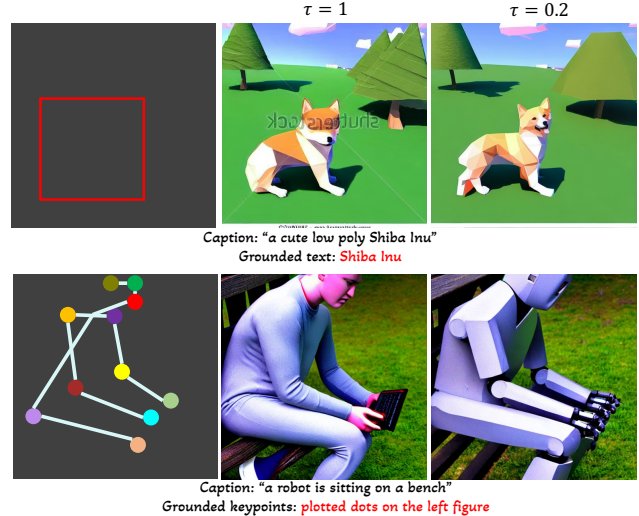


Figure 11. **Scheduled Sampling.** It can improve visual or extend a model trained in one domain (*e.g.*, human) to the others. See text for details.

set to 1), which is probably due to the grounding training data being all natural images. We believe adding images from more diverse style distributions or further finetuning the model with highly aesthetic images could help alleviate this issue.

**Acknowledgement.** We thank Yufan Zhou and Huangjie Zheng for the empirical insights on the FID evaluation of diffusion models, Haotian Zhang for the guidance on large-scale grounding data, Ce Liu for the discussion on the potential applications of grounded generative image models. This work was supported in part by NSF CAREER IIS2150012, NASA 80NSSC21K0295, and Institute of Information & communications Technology Planning & Evaluation(IITP) grants funded by the Korea government(MSIT) (No. 2022-0-00871, Development of AI Autonomy and Knowledge Enhancement for AI Agent Collaboration) and (No. RS-2022-00187238, Development of Large Korean Language Model Technology for Efficient Pre-training).

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. *ArXiv*, abs/2204.14198, 2022. 2, 5, 7
- [2] Oron Ashual, Shelly Sheynin, Adam Polyak, Uriel Singer, Oran Gafni, Eliya Nachmani, and Yaniv Taigman. Knn-

- diffusion: Image generation via large-scale retrieval. *arXiv preprint arXiv:2204.02849*, 2022. 6
- [3] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu Liu. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *ArXiv*, abs/2211.01324, 2022. 2, 8
  - [4] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 2
  - [5] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *ArXiv*, abs/2004.10934, 2020. 6
  - [6] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. 2
  - [7] Wenhui Chen, Hexiang Hu, Chitwan Saharia, and William W Cohen. Re-imagen: Retrieval-augmented text-to-image generator. *arXiv preprint arXiv:2209.14491*, 2022. 6
  - [8] Yunjei Choi, Min-Je Choi, Mun Su Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8789–8797, 2018. 3
  - [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. 3
  - [10] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. *ArXiv*, abs/2105.05233, 2021. 2
  - [11] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, and Jie Tang. Cogview: Mastering text-to-image generation via transformers, 2021. 6
  - [12] Wanshu 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*, abs/2208.13753, 2022. 6
  - [13] 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. *ArXiv*, abs/2203.13131, 2022. 2, 6
  - [14] 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 *NIPS*, 2014. 1, 3
  - [15] Agrim Gupta, Piotr Dollár, and Ross B. Girshick. Lvis: A dataset for large vocabulary instance segmentation. *CVPR*, pages 5351–5359, 2019. 2, 8, 15
  - [16] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 2
  - [17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 2
  - [18] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask r-cnn. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, 2017. 10
  - [19] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CVPR*, pages 770–778, 2016. 3
  - [20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NIPS*, 2017. 6
  - [21] Jonathan Ho. Classifier-free diffusion guidance. *ArXiv*, abs/2207.12598, 2022. 4, 6, 15
  - [22] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9:1735–1780, 1997. 3
  - [23] Ray S Jackendoff. *Semantic structures*, volume 18. MIT press, 1992. 4
  - [24] Manuel Jahn, Robin Rombach, and Björn Ommer. High-resolution complex scene synthesis with transformers. *ArXiv*, abs/2105.06458, 2021. 2, 3
  - [25] Manuel Jahn, Robin Rombach, and Björn Ommer. High-resolution complex scene synthesis with transformers. *ArXiv*, abs/2105.06458, 2021. 7, 16
  - [26] Justin Johnson, Agrim Gupta, and Li Fei-Fei. Image generation from scene graphs. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1219–1228, 2018. 2
  - [27] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *CVPR*, pages 4396–4405, 2019. 1
  - [28] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *CVPR*, pages 4396–4405, 2019. 3
  - [29] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8107–8116, 2020. 1
  - [30] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015. 15
  - [31] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123:32–73, 2016. 8
  - [32] Chunyuan Li, Haotian Liu, Liunian Harold Li, Pengchuan Zhang, Jyoti Aneja, Jianwei Yang, Ping Jin, Houdong Hu, Zicheng Liu, Yong Jae Lee, and Jianfeng Gao. ELEVATER: A benchmark and toolkit for evaluating language-augmented visual models. In *NeurIPS Track on Datasets and Benchmarks*, 2022. 2

- [33] Junnan Li, Ramprasaath R Selvaraju, Akhilesh Deepak Gotmare, Shafiq Joty, Caiming Xiong, and Steven Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *arXiv preprint arXiv:2107.07651*, 2021. [2](#)
- [34] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao. Grounded language-image pre-training. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 10955–10965. IEEE, 2022. [2](#), [6](#), [8](#)
- [35] Wenbo Li, Pengchuan Zhang, Lei Zhang, Qiuyuan Huang, Xiaodong He, Siwei Lyu, and Jianfeng Gao. Object-driven text-to-image synthesis via adversarial training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12174–12182, 2019. [6](#)
- [36] Yuheng Li, Yijun Li, Jingwan Lu, Eli Shechtman, Yong Jae Lee, and Krishna Kumar Singh. Collaging class-specific gans for semantic image synthesis. *ICCV*, pages 14398–14407, 2021. [3](#)
- [37] Yuheng Li, Yijun Li, Jingwan Lu, Eli Shechtman, Yong Jae Lee, and Krishna Kumar Singh. Contrastive learning for diverse disentangled foreground generation. *ArXiv*, abs/2211.02707, 2022. [2](#)
- [38] Yuheng Li, Krishna Kumar Singh, Utkarsh Ojha, and Yong Jae Lee. Mixnmatch: Multifactor disentanglement and encoding for conditional image generation. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8036–8045, 2020. [3](#)
- [39] Z. Li, Jingyu Wu, Immanuel Koh, Yongchuan Tang, and Lingyun Sun. Image synthesis from layout with locality-aware mask adaption. *ICCV*, pages 13799–13808, 2021. [2](#), [3](#)
- [40] Z. Li, Jingyu Wu, Immanuel Koh, Yongchuan Tang, and Lingyun Sun. Image synthesis from layout with locality-aware mask adaption. *ICCV*, pages 13799–13808, 2021. [6](#), [7](#), [8](#), [15](#), [16](#)
- [41] 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 *ECCV*, 2014. [2](#), [4](#), [6](#), [16](#)
- [42] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. [4](#), [5](#), [15](#)
- [43] 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. In *ICML*, 2022. [2](#), [11](#)
- [44] Vicente Ordonez, Girish Kulkarni, and Tamara L. Berg. Im2text: Describing images using 1 million captioned photographs. In *NIPS*, 2011. [8](#)
- [45] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. *CVPR*, pages 2332–2341, 2019. [1](#), [3](#)
- [46] Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. *CVPR*, pages 2536–2544, 2016. [1](#), [2](#)
- [47] Bryan A. Plummer, Liwei Wang, Christopher M. Cervantes, Juan C. Caicedo, J. Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. *International Journal of Computer Vision*, 123:74–93, 2015. [8](#)
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. [2](#), [3](#)
- [49] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *ArXiv*, abs/2204.06125, 2022. [1](#), [2](#), [6](#)
- [50] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8821–8831. PMLR, 18–24 Jul 2021. [1](#), [2](#)
- [51] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *CVPR*, pages 10674–10685, 2022. [2](#), [3](#), [6](#), [8](#), [15](#)
- [52] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, volume 9351 of *LNCS*, pages 234–241. Springer, 2015. (available on arXiv:1505.04597 [cs.CV]). [3](#), [15](#)
- [53] Chitwan Saharia, William Chan, Huiwen Chang, Chris A. Lee, Jonathan Ho, Tim Salimans, David J. Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. *ACM SIGGRAPH 2022 Conference Proceedings*, 2022. [2](#), [3](#)
- [54] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, Seyedeh Sara Mahdavi, Raphael Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *ArXiv*, abs/2205.11487, 2022. [1](#), [2](#), [6](#)
- [55] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. LAION-400M: open dataset of clip-filtered 400 million image-text pairs. *CoRR*, abs/2111.02114, 2021. [6](#)
- [56] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. *ICCV*, pages 8429–8438, 2019. [8](#)
- [57] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, 2018. [8](#)



- [58] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9240–9249, 2020. 1
- [59] Wei Sun and Tianfu Wu. Image synthesis from reconfigurable layout and style. *ICCV*, pages 10530–10539, 2019. 2, 3
- [60] Wei Sun and Tianfu Wu. Learning layout and style reconfigurable gans for controllable image synthesis. *TPAMI*, 44:5070–5087, 2022. 2, 3, 7, 16
- [61] Tristan Sylvain, Pengchuan Zhang, Yoshua Bengio, R. Devon Hjelm, and Shikhar Sharma. Object-centric image generation from layouts. *ArXiv*, abs/2003.07449, 2021. 2
- [62] Tristan Sylvain, Pengchuan Zhang, Yoshua Bengio, R. Devon Hjelm, and Shikhar Sharma. Object-centric image generation from layouts. *ArXiv*, abs/2003.07449, 2021. 7, 16
- [63] Ming Tao, Hao Tang, Songsong Wu, N. Sebe, Fei Wu, and Xiaoyuan Jing. Df-gan: Deep fusion generative adversarial networks for text-to-image synthesis. *ArXiv*, abs/2008.05865, 2020. 3, 6
- [64] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 3
- [65] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *ArXiv*, abs/1706.03762, 2017. 2, 3
- [66] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8798–8807, 2018. 8, 9
- [67] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. Nüwa: Visual synthesis pre-training for neural visual world creation. In *European Conference on Computer Vision*, 2022. 2, 6
- [68] 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. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1316–1324, 2018. 3, 6
- [69] Zuopeng Yang, Daqing Liu, Chaoyue Wang, J. Yang, and Dacheng Tao. Modeling image composition for complex scene generation. *CVPR*, pages 7754–7763, 2022. 2, 7, 16
- [70] Zuopeng Yang, Daqing Liu, Chaoyue Wang, J. Yang, and Dacheng Tao. Modeling image composition for complex scene generation. *CVPR*, pages 7754–7763, 2022. 2, 3
- [71] Zhengyuan Yang, Jianfeng Wang, Zhe Gan, Linjie Li, Kevin Lin, Chenfei Wu, Nan Duan, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Reco: Region-controlled text-to-image generation. *ArXiv*, abs/2211.15518, 2022. 3
- [72] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Benton C. Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation. *ArXiv*, abs/2206.10789, 2022. 1, 2, 6
- [73] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021. 2
- [74] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14393–14402, 2021. 2
- [75] Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, and Yinfei Yang. Cross-modal contrastive learning for text-to-image generation, 2021. 6
- [76] Bo Zhao, Lili Meng, Weidong Yin, and Leonid Sigal. Image generation from layout. *CVPR*, pages 8576–8585, 2019. 2, 3
- [77] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Llion Harold Li, Luwei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. RegionCLIP: Region-based language-image pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16793–16803, 2022. 2
- [78] Yufan Zhou, Chunyuan Li, Changyou Chen, Jianfeng Gao, and Jinhui Xu. Lafite2: Few-shot text-to-image generation. *arXiv preprint arXiv:2210.14124*, 2022. 3, 6
- [79] Yufan Zhou, Bingchen Liu, Yizhe Zhu, Xiao Yang, Changyou Chen, and Jinhui Xu. Shifted diffusion for text-to-image generation. *arXiv preprint arXiv:2211.15388*, 2022. 2
- [80] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. Towards language-free training for text-to-image generation. *CVPR*, 2022. 1, 6
- [81] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017. 3

## Appendix

In this supplemental material, we first provide more implementation and training details, and then present more results and discussions.

### A. Implementation and training details

We use the Stable Diffusion model [51] as the example to illustrate our implementation details.

**Box Grounding Tokens with Text.** Each grounded text is first fed into the text encoder to get the text embedding (e.g., 768 dimension of the CLIP text embedding in Stable Diffusion). Since the Stable Diffusion uses features of 77 text tokens outputted from the transformer backbone, thus we choose “EOS” token feature at this layer as our grounded text embedding. This is because in the CLIP training, this “EOS” token feature is chosen and applied a linear transform (one FC layer) to compare with visual feature, thus this token feature should contain whole information about the input text description. We also tried to directly use CLIP text embedding (after linear projection), however, we notice slow convergence empirically probably due to unaligned space between the grounded text embedding and the caption embeddings. Following NeRF [42], we encode bounding box coordinates with the Fourier embedding with output dimension 64. As stated in the Eq (5) in the main paper, we first concatenate these two features and feed them into a multi-layer perceptron. The MLP consists of three hidden layers with hidden dimension 512, the output grounding token dimension is set to be the same as the text embedding dimension (e.g., 768 in the Stable Diffusion case). We set the maximum number of grounding tokens to be 30 in the bounding box case.

**Box Grounding Tokens with Image.** We use the similar way to get the grounding token for an image. We use the CLIP image encoder (ViT-L-14 is used for the Stable Diffusion) to get an image embedding. We denote the CLIP training objective as maximizing  $(\mathbf{P}_t \mathbf{h}_t)^\top (\mathbf{P}_i \mathbf{h}_i)$  (we omit normalization), where  $\mathbf{h}_t$  is “EOS” token embedding from the text encoder,  $\mathbf{h}_i$  is “CLS” token embedding from the image encoder, and  $\mathbf{P}_t$  and  $\mathbf{P}_i$  are linear transformation for text and image embedding, respectively. Since  $\mathbf{h}_t$  is the text feature space used for grounded text features, to ease our training, we choose to project image features into the text feature space via  $\mathbf{P}_t^\top \mathbf{P}_i \mathbf{h}_i$ , and normalized it to 28.7, which is average norm of  $\mathbf{h}_t$  we empirically found. We also set the maximum number of grounding tokens to be 30. Thus, 60 tokens in total if one keep both image and text as representations for a grounded entity.

**Keypoint Grounding Tokens.** The grounding token for keypoint annotations is processed in the same way, ex-

cept that we also learn  $N$  person token embedding vectors  $\{\mathbf{p}_1, \dots, \mathbf{p}_N\}$  to semantically link keypoints belonging to the same person. This is to deal with the situation in which there are multiple people in the same image that we want to generate, so that the model knows which keypoint corresponds to which person. Each keypoint semantic embedding  $f_{\text{text}}(e)$  is processed by using the text encoder, for example, we forward the text: “left eye” into the encoder to get its semantic embedding; the dimension of each person token is set the same as text embedding dimension. The grounding token is calculated by:

$$\mathbf{h}^e = \text{MLP}(f_{\text{text}}(e) + \mathbf{p}_j, \text{Fourier}(\mathbf{l})) \quad (11)$$

where  $\mathbf{l}$  is the  $x, y$  location of each keypoint and  $\mathbf{p}_j$  is the person token for the  $j$ ’th person. In practice, we set  $N$  as 10, which is the maximum number of persons allowed to be generated in each image. Thus, we have 170 tokens in the COCO dataset (i.e.,  $10 \times 17$ ; 17 keypoint annotations for each person).

**Gated Self-Attention Layers.** Our inserted self-attention layer is the same as the original diffusion model self-attention layer at each Transformer block, except that we add one linear projection layer which converts the grounding token into the same dimension as the visual token. For example, in the first layer of the down branch of the UNet [52], the projection layer converts grounding token of dimension 768 into 320 (which is the image feature dimension at this layer), and visual tokens are concatenated with the grounding tokens as the input to the gated attention layer as illustrated in Figure 3.

**Training Details.** For all COCO related experiments (Sec 5.1), we train LDM with batch size 64 using 16 V100 GPUs for 100k iterations. In the scaling up training data experiment (in Sec. 5.2 of the main paper), we train for 400k iterations for LDM, but 500K iterations with batch size of 32 for the Stable diffusion model. For all training, we use learning rate of  $5e-5$  with Adam [30], and use warm-up for the first 10k iterations. We randomly drop caption and grounding tokens with 10% probability for classifier-free guidance [21].

### B. Additional quantitative results

In this section, we show more studies with our pretrained model using our largest data (GoldG, O365, CC3M, SBU). We had reported this model’s zero-shot performance on LVIS [15] in the main paper Table 3. Here we finetune this model on LVIS, and report its GLIP-score in Table 6. Clearly, after finetuning, we show much more accurate generation results, surpassing the supervised baseline LAMA [40] by a large margin.

Similarly, we also test this model’s zero-shot performance on the COCO2017 val-set, and its finetuning results are in

Model	Pre-training data	Training data	FID	AP	AP <sub>r</sub>	AP <sub>c</sub>	AP <sub>f</sub>
LAMA [40]	–	LVIS	151.96	2.0	0.9	1.3	3.2
GLIGEN-LDM	COCO2014CD	–	22.17	6.4	5.8	5.8	7.4
GLIGEN-LDM	COCO2014D	–	31.31	4.4	2.3	3.3	6.5
GLIGEN-LDM	COCO2014G	–	13.48	6.0	4.4	6.1	6.6
GLIGEN-LDM	GoldG,O365	–	8.45	10.6	5.8	9.6	13.8
GLIGEN-LDM	GoldG,O365,SBU,CC3M	–	10.28	11.1	9.0	9.8	13.4
GLIGEN-LDM	GoldG,O365,SBU,CC3M	LVIS	<b>6.25</b>	<b>14.9</b>	<b>10.1</b>	<b>12.8</b>	<b>19.3</b>
Upper-bound	–	–	–	25.2	19.0	22.2	31.2

Table 6. GLIP-score on LVIS validation set. Upper-bound is provided by running GLIP on real images scaled to  $256 \times 256$ .

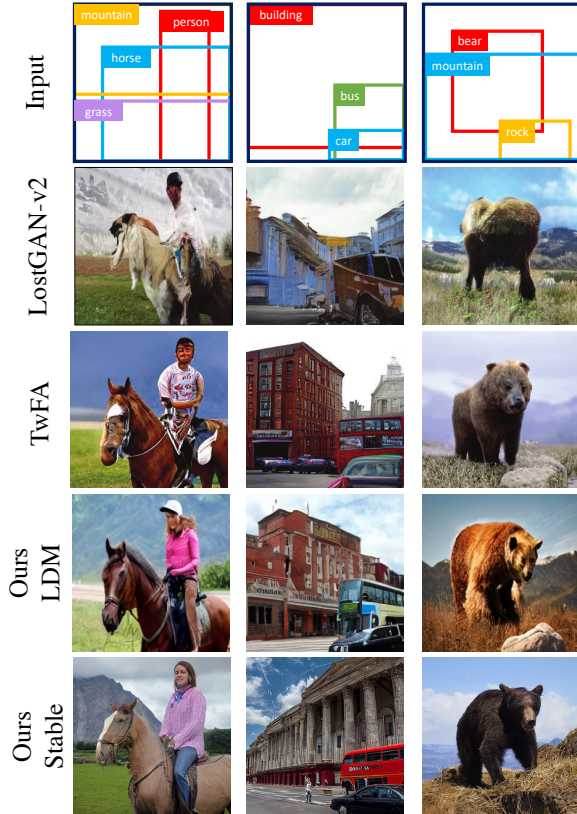


Figure 12. Layout2img comparison. Our model generates better quality images, especially when using stable diffusion. Baseline images are all copied from TwFA [69]

Table 7. The results show the benefits of pretraining which can largely improve layout correspondence performance.

Model	FID	YOLO score		
		AP	AP <sub>50</sub>	AP <sub>75</sub>
LostGAN-V2 [60]	42.55	9.1	15.3	9.8
OCGAN [62]	41.65	–	–	–
HCSS [25]	33.68	–	–	–
LAMA [40]	31.12	13.40	19.70	14.90
TwFA [69]	22.15	–	28.20	20.12
GLIGEN-LDM	<b>21.04</b>	22.4	36.5	24.1
<i>After pretrain on GoldG,O365,SBU,CC3M</i>				
GLIGEN-LDM (zero-shot)	27.03	19.1	30.5	20.8
GLIGEN-LDM (finetuned)	21.58	<b>30.8</b>	<b>42.3</b>	<b>35.3</b>

Table 7. Image quality and correspondence to layout are compared with baselines on COCO2017 val-set.

Lastly, we show more grounded text2img results with bounding boxes in Figure 13 and with images or keypoints in Figure 14. Note that our keypoint model only uses keypoint annotations from COCO [41] which is not linked with person identity, but it can successfully utilize and combine the knowledge learned in the text2img training stage to control keypoints of a specific person. Out of curiosity, we also tested whether the keypoint grounding information learned on humans can be transferred to other non-humanoid categories such as cat or lamp for keypoint grounded generation, but we find that our model struggles in such cases even with scheduled sampling. Compared to bounding boxes, which only specify a coarse location and size of an object in the image and thus can be shared across all object categories, keypoints (i.e., object parts) are not always shareable across different categories. Thus, while keypoints enable more fine-grained control than boxes, they are less generalizable.

### C. More qualitative results and discussion

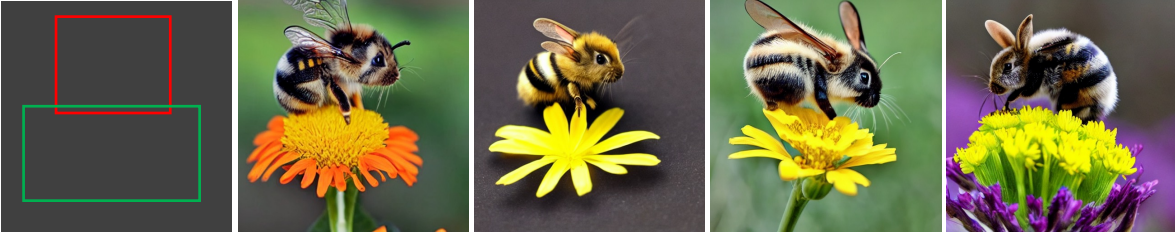
We show qualitative comparisons with layout2img baselines in Figure 12, which complements the results in Sec 5.1 of the main paper. The results show that our model has comparable image quality when built upon LDM, but has more visual appeal and details when built upon the Stable Diffusion model.





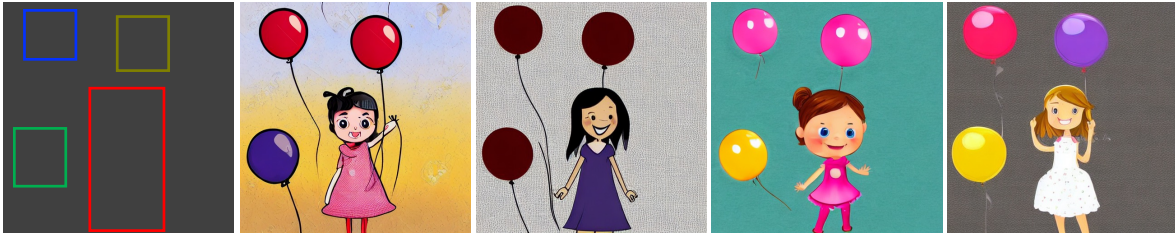
Caption: "Space view of a planet and its sun"

Grounded text: planet, sun



Caption: "a photo of a hybrid between a bee and a rabbit"

Grounded text: hybrid between a bee and a rabbit, flower



Caption: "cartoon sketch of a little girl with a smile and balloons, old style, detailed, elegant, intricate"

Grounded text: girl with a smile, balloon, balloon, balloon



Caption: "Walter White in GTA v"

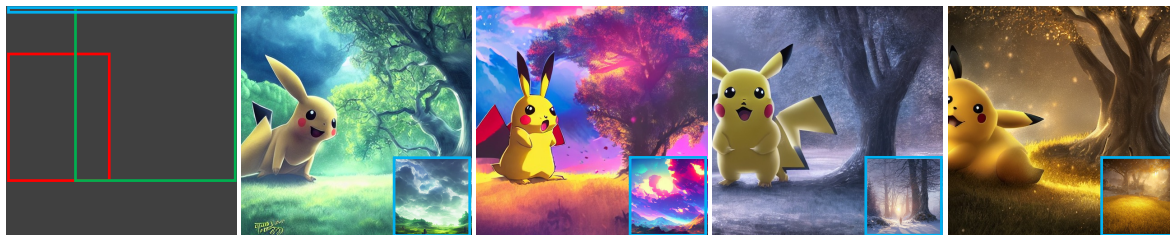
Grounded text: Walter White, car, bulldog



Caption: "two pirate ships on the ocean in minecraft"

Grounded text: a pirate ship, a pirate ship

Figure 13. Bounding box grounded text2image generation. Our model can ground noun entities in the caption for controllable image generation



Caption: "Pikachu is under a tree, digital art"  
 Grounded text: **Pikachu**, **tree**; Grounded image: **blue inset**



Caption: "Steve Jobs is working with his laptop"  
 Grounded keypoints: **plotted dots on the left**



Caption: "Barack Obama is sitting at a desk"  
 Grounded keypoints: **plotted dots on the left**



Figure 14. Top row: text grounded generation with style transfer. Second row: image grounded inpainting. Bottom rows: Keypoint grounded text2image generation. We did not use keypoint annotation associated with any identity during training.