

# Training-Free Layout-to-Image Generation with Marginal Attention Constraints

Huancheng Chen

University of Texas at Austin

huanchengch@utexas.edu

Jingtao Li

Sony AI

jingtao.li@sony.com

Weiming Zhuang

Sony AI

weiming.zhuang@sony.com

Haris Vikalo

University of Texas at Austin

hvikalo@ece.utexas.edu

Lingjuan Lyu\*

Sony AI

lingjuan.lv@sony.com

## Abstract

Recently, many text-to-image diffusion models excel at generating high-resolution images from text but struggle with precise control over spatial composition and object counting. To address these challenges, prior works developed layout-to-image (L2I) approaches that incorporate layout instructions into text-to-image models. However, existing L2I methods typically require fine-tuning of pre-trained parameters or training additional control modules for the diffusion models. In this work, we propose a training-free L2I approach, MAC (Marginal Attention Constrained Generation), which eliminates the need for additional modules or fine-tuning. Specifically, we use text-visual cross-attention feature maps to quantify inconsistencies between the layout of the generated images and the provided instructions, and then compute loss functions to optimize latent features during the diffusion reverse process. To enhance spatial controllability and mitigate semantic failures in complex layout instructions, we leverage pixel-to-pixel correlations in the self-attention feature maps to align cross-attention maps and combine three loss functions constrained by boundary attention to update latent features. Comprehensive experimental results on both L2I and non-L2I pretrained diffusion models demonstrate that our method outperforms existing training-free L2I techniques both quantitatively and qualitatively in terms of image composition on the DrawBench and HRS benchmarks.

## 1. Introduction

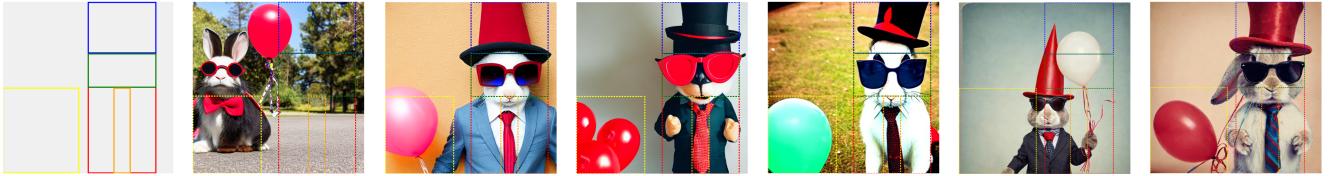
In recent years, text-to-image (T2I) diffusion models such as DALL-E [21], Imagen [24], and Stable Diffusion [22] have demonstrated an impressive ability to generate high-resolution images from textual input. These models derive their ability to unify textual and visual latent spaces from supervised training on large-scale datasets [20, 25] of text-

image pairs sourced from the internet. While such generative models achieve remarkable success in various downstream tasks, textual input alone cannot accurately specify spatial composition or the precise location of different concepts within generated images.

To address this deficiency, numerous studies [4, 13, 31, 32, 34] have explored layout-to-image (L2I) generation approaches that enable the localization of concept positions in prompts using various forms of layout instructions such as semantic masks, bounding boxes, or sketches. By incorporating additional layout information, text-to-image (T2I) models can achieve more precise spatial controllability and generate datasets with specific ground-truth labels for data augmentation in supervised training [13, 31]. However, these L2I approaches require adapting the pretrained T2I models with additional control modules that must be fine-tuned on large datasets containing paired images and layout annotations such as COCO [14].

The significant data and computational requirements for model fine-tuning restrict the use of such approaches in many data-scarce or resource-constrained scenarios. To this end, a series of methods for layout-to-image generation in a zero-shot manner, i.e., requiring no additional supervised training, have been proposed [5, 18, 29, 30, 33]. These techniques leverage the cross-attention maps extracted from U-Net [23] in the diffusion model to quantify the discrepancy between a synthetic image, sampled from the initialized latent feature  $z_t$  based on the input text prompt, and the target layout instructions; the subsequent iterative update of the latent feature  $z_t$  helps reduce the discrepancy. However, it has been observed that cross-attention maps typically exhibit high scores in the central regions of the concepts while assigning negligible scores to their edges [18, 29, 33]. Consequently, the generated concepts often appear larger than or misaligned with the specified bounding boxes, which diminishes the accuracy of layout control and leads to unsatisfactory image generation. Additionally, existing ap-

Prompt: "A **rabbit** wearing a red magician **hat**, **sunglasses**, and a **tie**, a bundle of **balloon**."



Prompt: "As the **aurora** lights up the sky, a herd of **reindeer** leisurely wanders on the grassy **meadow**, admiring the breathtaking view, a serene **lake** quietly reflects the magnificent display, and in the distance, a snow-capped mountain stands majestically."



Figure 1. An illustration of training-free layout-to-image (L2I) generation using various diffusion models. Text prompts and layout information, where specific concepts are localized within corresponding bounding boxes, are provided as input to the pretrained diffusion models to generate images that align with the instructions. Our method, MAC, is capable of guiding non-L2I diffusion models such as SD [22] and SDXL [19] to generate images based on layout instructions, while also enhancing L2I models such as Glichen [13] to achieve improved spatial control.

proaches often face overlapping cross-attention maps when the bounding boxes for multiple objects within the same concept are closely spaced with minimal gaps, leading to inaccuracies in object counts relative to the input prompt.

To improve spatial controllability and address semantic failures, we propose MAC, which stands for Marginal Attention Constrained Generation, a novel training-free L2I approach. Specifically, we introduce two key design principles: (1) leveraging pixel-to-pixel correlations from visual self-attention maps to align the coarse-grained cross-attention maps, and (2) proposing a boundary attention constraint to address challenges of size misalignment and inaccurate counting. Self-attention maps capture pixel-to-pixel correlations in the visual features, which can be used to filter noisy cross-attention maps and enhance the edges of concepts with low attention scores. Meanwhile, the boundary attention constraint ensures that cross-attention for each single object remains within the boundaries of the boxes, while promoting the separation of the cross-attention maps of multiple objects within the same concept. Our comprehensive experimental results demonstrate that the proposed method outperforms existing training-free L2I methods both quantitatively and qualitatively in terms of image compositions (*spatial relationship, size, color, object counting*) on the **DrawBench** [24] and **HRS** [2] benchmarks. The main contribution of this paper can be summarized as follows:

- We investigate semantic failures in training-free L2I generation with complex layout inputs, focusing on the issue

of overlapping cross-attention maps, which leads to inaccurate object counting.

- We propose a novel method, MAC, that advances the L2I generation performance by incorporating self-attention enhancement to filter noisy cross-attention maps, and boundary attention constraints to prevent overlapping cross-attention maps.
- We perform comprehensive experiments comparing our method with the existing L2I approaches. The quantitative and qualitative experimental results demonstrate that our method achieves state-of-the-art performance compared to existing L2I techniques.

## 2. Related Work

### 2.1. Text-To-Image Diffusion Models

Recently, diffusion models [9, 26, 27] have achieved remarkable success in image generation, demonstrating greater stability and controllability compared to GANs [7]. Essentially, a diffusion model learns the dynamics of mapping natural image-like data to a prior (e.g., Gaussian) distribution. The seminal work LDM [22] facilitates high-resolution image synthesis by performing forward and reverse processes in the latent space rather than the pixel space, thereby addressing issues of low inference speed and high training costs. In this framework, a U-net [23] is trained as the denoiser  $\epsilon_\theta$  to approximate the noise  $\epsilon_t$  based on the perturbed latent variable  $\mathbf{z}_t$  and the condition  $\mathbf{c}$  with

the loss

$$L(\theta) = \mathbb{E}_{\mathbf{z}_0, \epsilon \sim \mathcal{N}(0, 1), t} \left[ \|\epsilon_t - \epsilon_\theta(\mathbf{z}_t, \mathbf{c}, t)\|^2 \right]. \quad (1)$$

Building on the foundations of conditional generation with diffusion models, DALL-E [21] advances text-to-image (T2I) generation by utilizing the powerful image-language encoder CLIP [20] to transform text prompts into condition embeddings  $\mathbf{c}$  for guiding the reverse process of diffusion models. Imagen [24] and Stable Diffusion [22] extend this work by enhancing the diffusion sampling technique, thus enabling more photorealistic and detailed image generation. However, relying solely on text prompts limits the ability to precisely control spatial composition, necessitating additional input in the form of layout instructions.

## 2.2. Layout-To-Image Generation

In addition to text prompts, layout-to-image (L2I) models require auxiliary layout instructions as input, *i.e.*, semantic masks or bounding boxes. Several approaches aim to map layout instructions into the condition embedding space of diffusion models through supervised fine-tuning on datasets consisting of paired layout and image data. These approaches either integrate bounding boxes into the text prompts [1, 4, 31, 34] and fine-tune the text encoder, or incorporate additional modules or adapters alongside the text encoder to enhance the understanding of layout instructions [13, 32]. Another line of training-free approaches aims to address the computational and data challenges associated with supervised L2I methods. Inspired by the connection between spatial layouts of generated images and the cross-attention maps observed in [8], DenseDiff [12] guides the placement of objects in targeted positions by manipulating the cross-attention maps of the text encoder. Follow-up studies such as A&E [3] and BoxDiff [30] propose optimizing the latent variable  $z_t$  by maximizing the supremum of cross-attention scores located in specific regions during the reverse process of diffusion models. A&R [18] incorporates self-attention maps into the objective function to penalize objects that appear outside the designated boxes.

In contrast to previous studies that utilized the supremum, layout-guidance [34] uses the average of cross-attention scores in the objective function to update the latent variable  $z_t$ . Due to the typically coarse-grained and noisy nature of raw cross-attention maps, R&B [29] employs the Sobel Operator [10] to detect edges within these maps and select candidate boxes that encompass all edges for loss calculation. To address the additional computation and slower inference speed introduced by edge detection, concurrent work LoCo [33] directly leverages start-of-text tokens (SoT), end-of-text tokens (EoT), and self-attention maps to enhance cross-attention maps. B2B [28] incorporates objection generation and attribute binding to improve the controllability. Concurrent work [6] explores the

attention leakage problem and uses bounded attention for reducing semantic leakage. However, none of the aforementioned approaches incorporates marginal attention constraints, leading to failures in object counting and precise spatial controllability.

## 3. Methodology

### 3.1. Training-Free Layout-To-Image Generation

Training-free layout-to-image generation aims to generate images based on input text prompts and corresponding layout instructions using a pretrained diffusion model, without the need for additional parameter training or extra modules. Let us consider a text prompt  $\mathbf{p} = \{\mathbf{p}_1, \dots, \mathbf{p}_n\}$  and a set of bounding boxes  $\mathcal{B} = \{\mathcal{B}_i \text{ for } \forall i \in \mathcal{I}, \text{ where } \mathcal{I} \subset [n]\}$ ; here  $\mathcal{B}_i = \{\mathbf{b}_1^{(i)}, \dots, \mathbf{b}_{N_i}^{(i)}\}$  denotes the corresponding bounding boxes for phrase  $\mathbf{p}_i$  consisting of  $N_i$  pairs of top-left and bottom-right points  $(x_1, y_1, x_2, y_2)$  that specify locations of  $N_i$  objects described by  $\mathbf{p}_i$ . The generated images should closely align with the text prompts while being consistent with the layout instructions defined by the set of bounding boxes  $\mathcal{B}$ . An illustration of training-free layout-to-image generation using bounding boxes as layout instructions is given in Figure 1.

### 3.2. Cross-Attention and Self-Attention Maps

In the T2I diffusion models, text prompt  $\mathbf{p}$  is encoded via text tokens  $\mathbf{e} = f_{\text{CLIP}}(\mathbf{p}) \in \mathbb{R}^{n \times d_e}$  using a pretrained CLIP encoder. These tokens are then input into the cross-attention layers, where they are fused with the visual embedding characterized by  $\mathbf{z}_t$ . Specifically, in layer  $l$ , the visual and text embeddings are projected into the *query*  $\mathbf{Q}_z^l \in \mathbb{R}^{hw \times d}$  and *key*  $\mathbf{K}_e^l \in \mathbb{R}^{n \times d}$ , respectively, allowing the cross-attention maps to be computed as

$$\mathbf{A}_c^l = \text{softmax} \left( \frac{\mathbf{Q}_z^l (\mathbf{K}_e^l)^\top}{\sqrt{d}} \right) \in [0, 1]^{hw \times n}, \quad (2)$$

where  $n$  denotes the length of the text prompt (including SoT and EoT);  $d_e$  and  $d$  denote the dimensions of the text and visual embeddings, respectively; and  $h$  and  $w$  represent the height and width of the visual feature maps, respectively. Similarly, the image embeddings are projected into the *key* matrix  $\mathbf{K}_z^l \in \mathbb{R}^{hw \times d}$ , enabling computation of self-attention maps

$$\mathbf{A}_s^l = \text{softmax} \left( \frac{\mathbf{Q}_z^l (\mathbf{K}_z^l)^\top}{\sqrt{d}} \right) \in [0, 1]^{hw \times hw} \quad (3)$$

which measures the pixel-to-pixel similarity in the visual features. These cross-attention and self-attention maps obtained from  $L$  different cross-attention layers are aggregated as

$$\mathbf{A}_c = \frac{1}{L} \sum_{l=1}^L \mathbf{A}_c^l, \quad \mathbf{A}_s = \frac{1}{L} \sum_{l=1}^L \mathbf{A}_s^l. \quad (4)$$

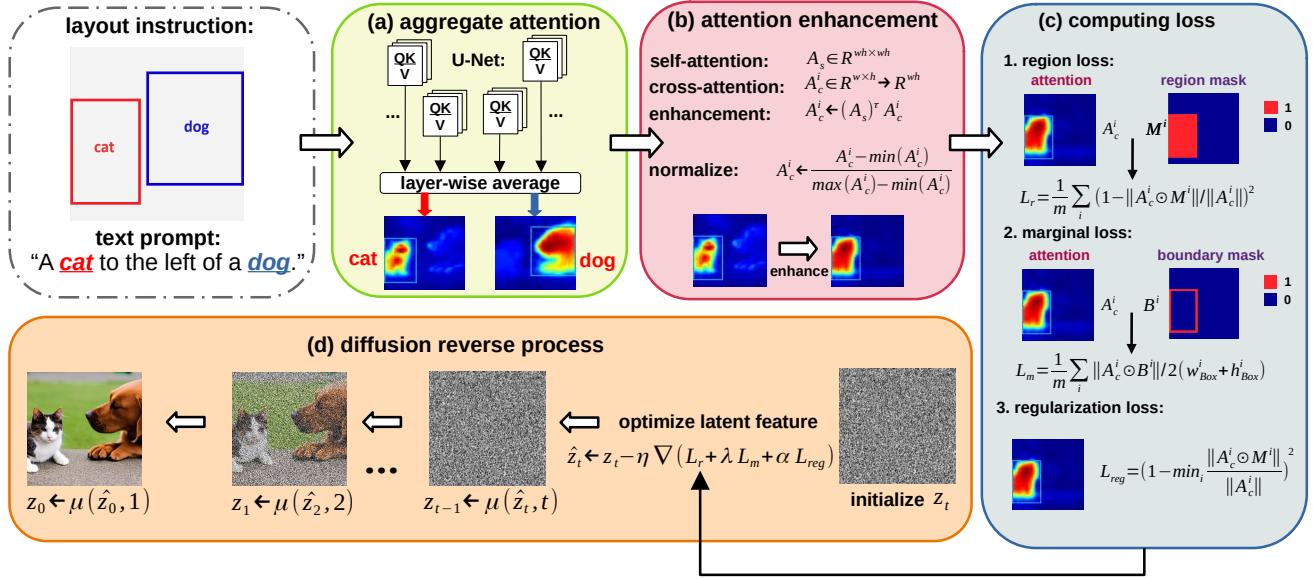


Figure 2. The overall framework of MAC for training-free L2I generation consists of four steps: (a) aggregating cross-attention and self-attention maps obtained from the language-vision fusion blocks by layer-wise averaging; (b) enhancing cross-attention maps with self-attention maps; (c) computing three losses and combining them with coefficients; (d) optimizing  $z_t$  by gradient descent.

### 3.3. Marginal Attention Constrained Guidance

While prior studies [5, 29] demonstrated effectiveness of backward guidance, a procedure where cross-attention scores are used to form the loss utilized for optimizing the latent feature  $z_t$  during the diffusion reverse process, the designed loss functions fail to address issues related to counting and semantic failures arising due to the coarse-grained nature of the cross-attention maps. In this section, we discuss the effect of attention enhancement and present three loss functions to address the problems encountered by previous schemes.

**Attention Enhancement.** As depicted in the pink block in Figure 2, the raw cross-attention map for “cat” is coarse-grained, with attention dispersed across multiple parts of the object. It also shows notable scores for “dog” in the cross-attention map for the “cat”, likely due to the high correlation between “cat” and “dog” within the pretrained model, which causes inconsistency in the loss computation. Inspired by the previous studies [16, 17] that aim to build synthetic datasets by generating grounded labels using cross-attention, we leverage the global information contained in self-attention maps to enhance the cross-attention maps,

$$\mathbf{A}_c^{(i)} = (\mathbf{A}_s)^{\tau} \cdot \mathbf{A}_c^{(i)} \in [0, 1]^{wh}, \quad (5)$$

where  $i$  denotes the index of the phrase  $\mathbf{p}_i$ , and  $\tau$  is a power coefficient for adjusting the magnitude of enhancement (set to 1 by default). For dimensional consistency, we normalize

and reshape  $\mathbf{A}_c^{(i)}$  according to

$$\mathbf{A}_c^{(i)} = \text{reshape} \left( \frac{\mathbf{A}_c^{(i)} - \min(\mathbf{A}_c^{(i)})}{\max(\mathbf{A}_c^{(i)}) - \min(\mathbf{A}_c^{(i)})} \right) \in [0, 1]^{w \times h}. \quad (6)$$

Essentially, the self-attention map  $\mathbf{A}_s$  captures the pairwise correlations between pixels, allowing salient parts in the raw cross-attention maps to propagate to the most relevant regions. Using these enhanced cross-attention maps, we compute three losses for improving the L2I generation.

**Region-Attention Loss.** Similar to the prior studies [5, 18, 29], the key idea of region-attention loss is to measure the proportion of cross-attention map  $\mathbf{A}_c^{(i)}$  outside bounding boxes  $\mathcal{B}_i \in \mathcal{B}$  and compute the average of these proportions for normalization as

$$\mathbf{M}^{(i)}[x, y] = \begin{cases} 1, & \text{if } (x, y) \text{ inside } \mathcal{B}_i \\ 0, & \text{if } (x, y) \text{ outside } \mathcal{B}_i \end{cases}, \quad (7)$$

$$L_r = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{I}} \left( 1 - \frac{\mathbf{A}_c^{(i)} \odot \mathbf{M}^{(i)}}{\sum_{x,y} \mathbf{A}_c^{(i)}[x, y]} \right)^2, \quad (8)$$

where  $\mathbf{A}_c^{(i)}[x, y]$  and  $\mathbf{M}^{(i)}[x, y]$  are the  $(x, y)$ -th element in  $\mathbf{A}_c^{(i)}$  and  $\mathbf{M}^{(i)}$ , respectively. However, relying solely on region-attention loss  $L_r$  cannot ensure high-quality generation in scenarios where multiple objects constrained by the same  $\mathcal{B}_i$  have overlapping cross-attention scores, resulting in incorrect counting in the generated image.

**Marginal-Attention Loss.** As previously mentioned, adjacent bounding boxes for multiple objects described by the same phrase often lead to incorrect counting in the L2I generation due to the interference between cross-attention scores. However, the region-attention loss may remain small as long as these cross-attention scores are inside the target bounding boxes. We propose a marginal-attention loss that aims to isolate adjacent cross-attention maps corresponding to different objects,

$$\mathbf{B}^{(i)}[x, y] = \begin{cases} 1, & \text{if } (x, y) \text{ on } \mathcal{B}_i \\ 0, & \text{if } (x, y) \text{ not on } \mathcal{B}_i \end{cases}, \quad (9)$$

$$L_m = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{I}} \frac{\mathbf{A}_c^{(i)} \odot \mathbf{B}^{(i)}}{2(w_{\text{Box}} + h_{\text{Box}})}, \quad (10)$$

where  $w_{\text{Box}} = \sum_{k=1}^{N_i} w_k^{(i)}$ ,  $h_{\text{Box}} = \sum_{k=1}^{N_i} h_k^{(i)}$ , and  $w_k^{(i)}$  and  $h_k^{(i)}$  denote the width and height of the bounding box  $\mathbf{b}_k^{(i)}$ , respectively. With the proposed marginal-attention loss, we can force  $z_t$  to generate cross-attention maps without overlap between multiple objects and thus help improve the accuracy of counting.

**Regularization Loss.** While the region-attention and marginal-attention losses generally help localize generated objects within the target bounding boxes, these losses may be minimal when an incorrect number of objects has cross-attention scores completely contained within the bounding boxes (i.e., experiencing no boundary crossing). To prevent bad generation in such scenarios, we propose a regularized loss that ensures the assigned objects do not have a void cross-attention map

$$L_{\text{reg}} = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{I}} \left( 1 - \min_k \frac{\mathbf{A}_c^{(i)} \odot \mathbf{M}^{(i,k)}}{\sum_{x,y} \mathbf{A}_c^{(i)}[x, y]} \right)^2, \quad (11)$$

where  $\mathbf{M}^{(i,k)}$  is a subset of  $\mathbf{M}^{(i)}$  indicating the location of the  $k$ -th object in  $\mathcal{B}_i$ .

**Latent Feature Optimization.** At each sampling time step in the reverse process,  $L_{\text{mac}}$  can be computed as the weighted summation of  $L_r$ ,  $L_m$  and  $L_{\text{reg}}$ ,

$$L_{\text{mac}} = L_r + \lambda L_m + \alpha L_{\text{reg}}, \quad (12)$$

where  $\lambda$  and  $\alpha$  control the intervention strength of  $L_m$  and  $L_{\text{reg}}$ , respectively. We compute and use the gradient of  $L_{\text{mac}}$  to update  $z_t$  as  $\hat{z}_t \leftarrow z_t - \eta \nabla L_{\text{mac}}$ , where parameter  $\eta$  controls the size of the updates. After  $t_{\text{optim}}$  iterations or an early stop that happens when  $L_{\text{mac}}$  becomes smaller than a predetermined threshold, the optimized latent feature  $\hat{z}_t$  is forwarded to the U-Net for predicting the latent feature  $z_{t-1}$  at the previous time step.

## 4. Experiments

### 4.1. Experimental Setup

**Datasets and base models.** Following the strategy in the prior studies [18, 29], we perform experiments on a subset of two widely used benchmarks **HRS** [2] and **Drawbench** [24]. Specifically, we use four tracks from the HRS benchmark – **spatial relationship**, **size**, **color**, and **object counting** – which include 1002, 501, 501, and 3000 text prompts, respectively. These prompts, along with the corresponding layout instructions generated by ChatGPT-4 in [18], are used to evaluate the performance of our method and baseline approaches both quantitatively and qualitatively. Similarly, 39 text prompts from Drawbench involving **spatial** and **counting** specifications are also used in the evaluation. We conduct experiments with the widely-used Stable Diffusion (SD) 1.5 [22] as the base model and perform an ablation study with SD-XL [19] and the fine-tuned L2I model, Gligen [13]. The sampling time step is set to 50, with the classifier-free guidance weight set to 7.5. Optimization is applied only to the latent feature  $z_t$  during the first 10 sampling steps, and  $\eta$  is set to 70. Unless specified otherwise, the hyperparameters  $\lambda$  and  $\alpha$  are both set to 0.5.

**Baselines and metrics.** We compare our proposed scheme, MAC, to six state-of-the-art methods: A&E [3], BoxDiff [30], Layout-Guidance [5], A&R [18], R&B [29], and LoCo [33]. To quantitatively assess our method against these baselines, we use the state-of-the-art object detector, Ground-DINO [15], to detect objects in the synthetic images and predict bounding boxes. These predicted boxes are then compared with the ground truth. We compute precision, recall, and F1 metrics to comprehensively evaluate MAC’s performance in object counting. By comparing the area and centroid of the predicted and ground-truth boxes, we compute the accuracy of object size and spatial relationships. Additionally, Ground-DINO predicts the color of detected objects, which we use to compute color accuracy. The details of the metric computations can be found in the Appendix B.

### 4.2. Visual Comparisons

In Figure 3, we present visual comparisons between MAC and several competing L2I schemes – Layout, A&R, and R&B – deployed on SD-XL [19]. To thoroughly evaluate the performance of these schemes, we use several complex inputs featuring multiple objects distributed in random positions, which makes it more challenging to generate images with precise object counts. As shown in the figure, MAC demonstrates significant improvement over the baselines in terms of (1) better alignment within the bounding boxes, and (2) more precise object counting. Although the prior works typically generate objects within the corresponding

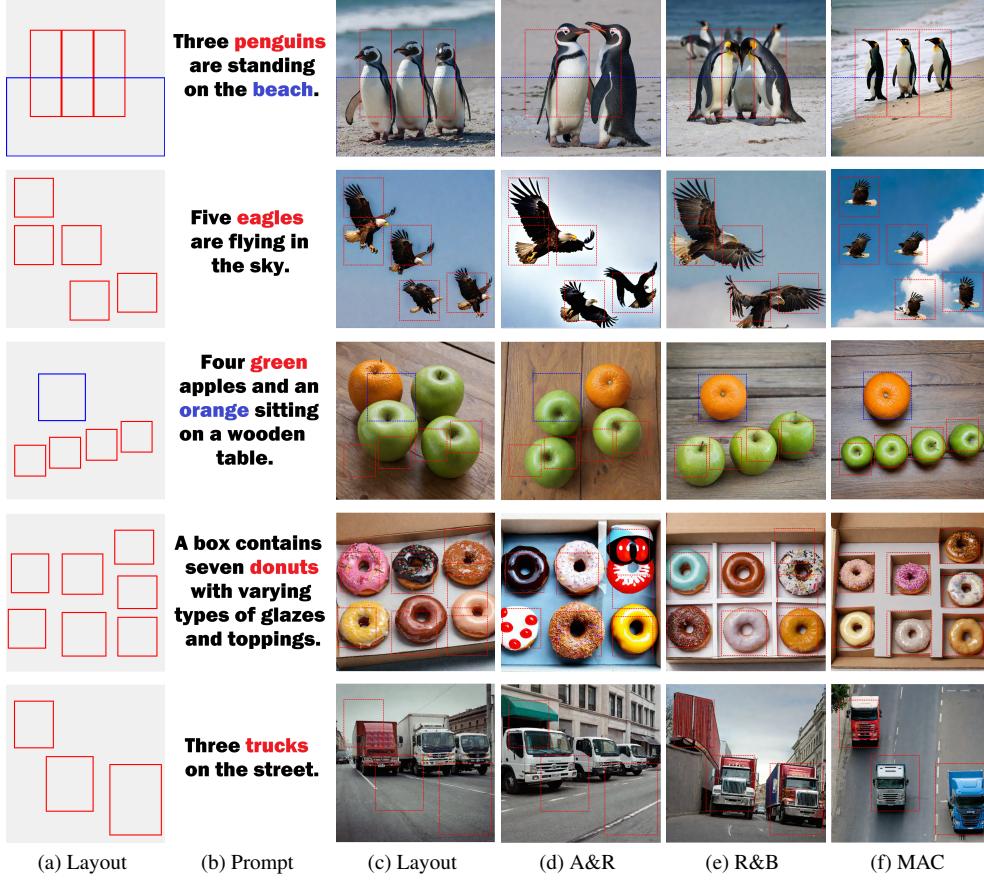


Figure 3. Visual comparison with different training-free L2I schemes based on Stable Diffusion XL [19]. The layout input and text prompt are shown in the first and second columns, respectively. The layout instructions are also annotated on the generated images with dashed boxes. Our method, MAC, significantly outperforms prior schemes in terms of spatial control and counting accuracy.

bounding boxes, most of these objects are only partially inside the boxes or even appear in the bounding boxes of other objects due to coarse-grained cross-attention maps, as we discussed earlier. For example, Layout [34] and A&R [18] fail to generate penguins (1st row) within the red bounding boxes and an orange (3rd row) within the blue bounding box. Additionally, none of these prior schemes can generate the correct number of objects according to the prompts and bounding boxes specifications.

To better understand the miscounting failures, we visualize the cross-attention maps obtained from these schemes as shown in Figure 4. Since the cross-attention maps for objects described by the same phrase in the prompt are not independent, existing training-free L2I schemes, which aim to encourage cross-attention within the corresponding boxes, often produce overlapping cross-attention maps. For instance, in the R&B scheme, the cross-attention scores are concentrated within three boxes, representing a single eagle, even though three eagles are meant to be depicted. However, marginal-attention loss proposed in MAC can isolate

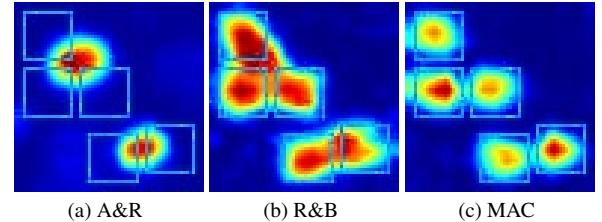


Figure 4. Visualization of cross-attention maps with the prompt: *Five eagles are flying in the sky*.

these cross-attention scores to ensure they depict independent objects in the generated images.

### 4.3. Quantitative Results

We validate the improvement in spatial controllability and object counting of the proposed MAC framework as compared to existing state-of-the-art training-free L2I approaches through quantitative experiments. As shown in Tables 2 and 1, MAC outperforms state-of-the-art baselines

Table 1. Quantitative comparisons with baseline schemes. The precision (%), recall (%) and F1 score (%) with respect to object counting are reported in the table.

Method	DrawBench			HRS-Bench			Inference (s)
	Precision	Recall	F1	Precision	Recall	F1	
SD1.5	76.51	70.58	71.35	71.86	52.19	58.31	14
+ A&E	75.52	66.17	76.27	73.10	54.79	60.47	42
+ BoxDiff	86.36	77.94	86.17	73.57	63.33	68.07	42
+ LoCo	86.04	54.41	66.66	85.61	59.55	71.23	39
+ Layout	85.93	80.88	83.33	87.25	63.10	73.24	<b>37</b>
+ A&R	86.00	63.23	72.88	87.93	56.69	68.94	45
+ R&B	87.50	82.35	84.84	85.73	68.53	76.17	84
+ MAC	<b>93.84</b>	<b>89.71</b>	<b>91.72</b>	<b>89.51</b>	<b>68.92</b>	<b>77.88</b>	41

Table 2. Quantitative comparisons with baseline schemes. The evaluated accuracy (%) with respect to spatial relationship, size relationship and color are reported in the table.

Benchmarks	Metrics	SD	A&E	BoxDiff	LoCo	Layout	A&R	R&B	MAC
DrawBench	Spatial	12.50	15.00	25.00	40.00	45.00	40.00	55.00	<b>60.00</b>
	Size	11.23	14.77	12.77	14.56	14.37	12.17	16.17	<b>16.96</b>
HRS-Bench	Color	13.01	18.27	34.69	37.88	31.18	30.15	36.36	<b>39.74</b>
	Spatial	10.80	17.56	19.76	37.42	33.73	25.94	47.80	<b>50.39</b>

on both DrawBench and HRS benchmarks. Notably, MAC demonstrates significant improvements in object counting across three metrics: precision, recall, and F1 score. On DrawBench, MAC outperforms the best baseline, R&B, by 6.34%, 7.36%, and 6.88% in precision, recall, and F1 score, respectively. On the larger HRS benchmark, MAC consistently achieves the best performance among all L2I schemes, with improvements of 3.78%, 0.39%, and 1.71% in precision, recall, and F1 score, respectively. Furthermore, MAC surpasses R&B in the spatial relationship metric on both DrawBench and HRS benchmark, with improvement of 5% and 2.59%. Additionally, MAC achieves the best performance in color and size metrics, though its advantage over R&B in these areas is relatively small.

Nevertheless, R&B requires applying Sobel operator [11] to detect the edge of cross-attention map for creating a minimum bounding rectangle (MBR), which is then used to compute region-aware loss. Therefore, R&B needs more computation in the process of optimizing latent feature  $z_t$ , which leads to longer inference time. As shown in Table 2, R&B needs more than double time on average for generating one image compared to MAC, even though R&B has comparable performance to MAC.

#### 4.4. Plug and Play with MAC

Although we mainly conduct experiments based on SD 1.5, MAC, as well as these prior training-free L2I schemes do not assume specific model architecture and can be plugged into arbitrary pre-trained Text-to-Image models using cross-

attention blocks. We further evaluate the performance of MAC by deploying it on another L2I model, GLIGEN [13], trained by supervised learning with labeled datasets. As shown in Table 3, all training-free L2I schemes demonstrate improved performance across nearly all metrics quantifying controllability. Compared to SD 1.5, GLIGEN has the ability to perceive input bounding boxes and incorporate layout information into the prior conditions for sampling the initial noise  $z_0$ . With the initial latent feature  $z_0$  generated by GLIGEN, the reverse sampler can synthesize an image with a layout that more closely aligns with the given instructions than when using  $z_0$  produced by SD 1.5. According to Table 3, MAC consistently delivers the best performance in object counting, spatial relationships, and color. While R&B outperforms MAC in the size metric, MAC achieves the second-best performance. In experiments with SD 1.5, plain SD 1.5 generates initial noise  $z_0$  corresponding to undesirable cross-attention maps that are minimally modified through optimization, resulting in limited improvement of MAC on the size metric. In contrast, GLIGEN generates improved initial noise  $z_0$  that facilitates MAC’s enhancement.

#### 4.5. Ablation Study

To investigate the effect of  $L_r$ ,  $L_b$  and  $L_{reg}$ , the components of MAC’s loss / objective function, we conduct additional experiments using SD 1.5, SD-XL, and GLIGEN under three settings: (1) only the region-attention loss  $L_r$  is used (setting  $R$ ); (2)  $L_r$  is combined with the marginal-

Table 3. Zero-shot schemes can be integrated into existing fully-supervised layout-to-image schemes, e.g., GLIGEN [13]. We compare the improvement of GLIGEN augmented with MAC vs. the competing schemes.

Method	DrawBench		HRS-Bench			
	Counting (F1)	Spatial	Counting (F1)	Spatial	Size	Color
GLIGEN	77.03	40.00	68.32	44.81	35.31	30.55
+ LoCo	76.61	40.00	71.89	49.48	38.33	37.97
+ Layout	80.59	55.00	81.85	50.68	29.64	35.05
+ A&R	76.27	55.00	77.66	52.17	30.72	38.74
+ R&B	87.59	<b>65.00</b>	79.10	53.88	39.72	40.59
+ MAC	<b>93.22</b>	<b>65.00</b>	<b>83.34</b>	<b>55.76</b>	<b>43.31</b>	<b>42.48</b>

Table 4. The results of experiments on MAC under three settings: (1)  $L_{\text{mac}}$  only includes the region-aware loss (setting  $R$ ); (2)  $L_{\text{mac}}$  includes both the region-aware loss and marginal-aware loss (setting  $R + M$ ); (3) using the complete loss proposed in Eq. 12 (setting  $R + M + \text{Reg}$ ).

Method	HRS-Bench			
	Counting	Spatial	Size	Color
SD1.5	71.35	10.80	11.23	13.01
+ $R$	73.48	30.57	14.15	30.67
+ $R + M$	76.03	43.61	15.63	32.53
+ $R + M + \text{Reg}$	<b>77.88</b>	<b>50.39</b>	<b>16.96</b>	<b>39.74</b>
SD-XL	76.99	34.23	19.96	25.20
+ $R$	76.71	33.34	30.34	27.42
+ $R + M$	78.09	38.82	27.34	28.42
+ $R + M + \text{Reg}$	<b>80.43</b>	<b>42.31</b>	<b>38.12</b>	<b>32.12</b>
GLIGEN	68.32	44.81	35.31	30.55
+ $R$	79.94	49.43	30.87	35.46
+ $R + M$	82.63	53.08	38.12	40.44
+ $R + M + \text{Reg}$	<b>83.34</b>	<b>55.76</b>	<b>43.31</b>	<b>42.28</b>

attention loss  $L_m$  (setting  $R + M$ ); (3) using the complete loss objective as described in Eq. 12 (setting  $R + M + \text{Reg}$ ).

When using only the region-attention loss, MAC synthesizes images similar to those produced by R&B, with the cross-attention map exhibiting the same overlapping issue discussed earlier – see the illustration in Figure 5(a). Adding the marginal-attention loss to the objective resolves the overlapping problem but may cause attention vanishing, as shown in Figure 5(b). This issue does not occur when there is only one bounding box per phrase, due to the constraint imposed by  $L_r$ . Moreover, the regularization loss  $L_{\text{reg}}$  addresses the attention vanishing problem, ensuring that each bounding box contains exactly one object in the generated images, as shown in Figure 5(c).

The quantitative results in Table 4, obtained in experiments on three distinct pretrained diffusion models, further demonstrate the improvement achieved by including  $L_m$  and  $L_{\text{reg}}$  in the objective function. However, the attention

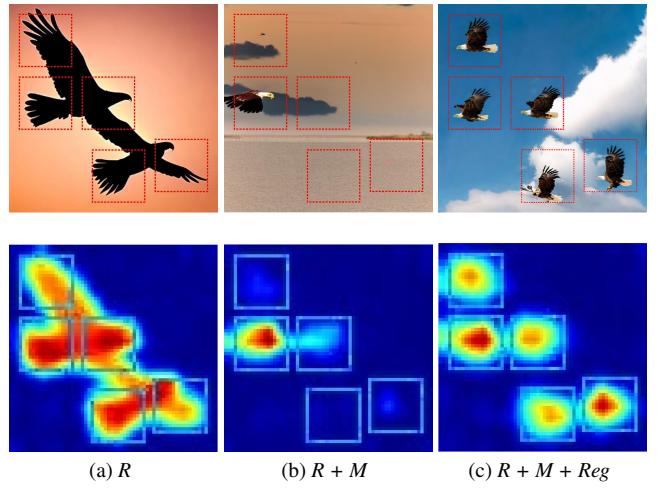


Figure 5. Visualization of cross-attention maps of targeting objects generated in three settings: (1)  $R$ ; (2)  $R + M$ ; (3)  $R + M + \text{Reg}$ .

vanishing problem typically arises in settings where multiple objects are assigned adjacent bounding boxes. The text prompts in HRS are relatively simple, so there is minimal difference between  $R + M$  and  $R + M + \text{Reg}$  settings.

## 5. Conclusion

We proposed MAC, a novel training-free layout-to-image (L2I) generation scheme utilizing pre-trained diffusion models. We identify two inherent issues that lead to semantic errors and incorrect object counts in generation: (1) coarse-grained cross-attention maps, and (2) overlapping cross-attention. To address these issues, we leverage self-attention maps to refine cross-attention maps and design an objective that combines three distinct loss functions. Comprehensive quantitative and visual results demonstrate that MAC outperforms state-of-the-art L2I schemes and consistently adapts well across varying model architectures. We envision our MAC as a milestone that can inspire the follow-up research in training-free layout-to-image generation.

## 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. 3
- [2] Eslam Mohamed Bakr, Pengzhan Sun, Xiaogian Shen, Faizan Farooq Khan, Li Erran Li, and Mohamed Elhoseiny. Hrs-bench: Holistic, reliable and scalable benchmark for text-to-image models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20041–20053, 2023. 2, 5
- [3] 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. 3, 5
- [4] Kai Chen, Enze Xie, Zhe Chen, Yibo Wang, Lanqing Hong, Zhenguo Li, and Dit-Yan Yeung. Geodiffusion: Text-prompted geometric control for object detection data generation. *arXiv preprint arXiv:2306.04607*, 2023. 1, 3
- [5] Minghao Chen, Iro Laina, and Andrea Vedaldi. Training-free layout control with cross-attention guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5343–5353, 2024. 1, 4, 5
- [6] Omer Dahary, Or Patashnik, Kfir Aberman, and Daniel Cohen-Or. Be yourself: Bounded attention for multi-subject text-to-image generation. In *European Conference on Computer Vision*, pages 432–448. Springer, 2024. 3
- [7] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. 2
- [8] 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. 3
- [9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2
- [10] Nick Kanopoulos, Nagesh Vasanthavada, and Robert L Baker. Design of an image edge detection filter using the sobel operator. *IEEE Journal of solid-state circuits*, 23(2): 358–367, 1988. 3
- [11] Nick Kanopoulos, Nagesh Vasanthavada, and Robert L Baker. Design of an image edge detection filter using the sobel operator. *IEEE Journal of solid-state circuits*, 23(2): 358–367, 1988. 7
- [12] Yunji Kim, Jiyoung Lee, Jin-Hwa Kim, Jung-Woo Ha, and Jun-Yan Zhu. Dense text-to-image generation with attention modulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7701–7711, 2023. 3
- [13] Yuheng Li, Haotian Liu, et al. Gligen: Open-set grounded text-to-image generation. *CVPR*, 2023. 1, 2, 3, 5, 7, 8
- [14] 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 *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 1
- [15] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 5, 11
- [16] Chaofan Ma, Yuhuan Yang, Chen Ju, Fei Zhang, Jinxiang Liu, Yu Wang, Ya Zhang, and Yanfeng Wang. Diffusionseg: Adapting diffusion towards unsupervised object discovery. *arXiv preprint arXiv:2303.09813*, 2023. 4
- [17] Quang Nguyen, Truong Vu, Anh Tran, and Khoi Nguyen. Dataset diffusion: Diffusion-based synthetic data generation for pixel-level semantic segmentation. *Advances in Neural Information Processing Systems*, 36, 2024. 4
- [18] Quynh Phung, Songwei Ge, and Jia-Bin Huang. Grounded text-to-image synthesis with attention refocusing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7932–7942, 2024. 1, 3, 4, 5, 6
- [19] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2, 5, 6
- [20] 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. 1, 3
- [21] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 1, 3
- [22] 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. 1, 2, 3, 5
- [23] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 1, 2
- [24] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022. 1, 2, 3, 5
- [25] 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. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 1
- [26] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 2, 11
- [27] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021. 2
- [28] Ashkan Taghipour, Morteza Ghahremani, Mohammed Ben-namoun, Aref Miri Rekavandi, Hamid Laga, and Farid Bousaid. Box it to bind it: Unified layout control and attribute binding in t2i diffusion models. *arXiv preprint arXiv:2402.17910*, 2024. 3
- [29] Jiayu Xiao, Liang Li, et al. R&b: Region and boundary aware zero-shot grounded text-to-image generation. *ICLR*, 2023. 1, 3, 4, 5
- [30] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7452–7461, 2023. 1, 3, 5
- [31] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14246–14255, 2023. 1, 3
- [32] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 1, 3
- [33] Peiangular Zhao, Han Li, Ruiyang Jin, and S Kevin Zhou. Loco: Locally constrained training-free layout-to-image synthesis. *arXiv preprint arXiv:2311.12342*, 2023. 1, 3, 5
- [34] Guangcong Zheng, Xianpan Zhou, Xuewei Li, Zhongang Qi, Ying Shan, and Xi Li. Layoutdiffusion: Controllable diffusion model for layout-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22490–22499, 2023. 1, 3, 6

## Appendix A: Experimental Details

In this section, we report the model architectures and hyper-parameters used in the experiments across all layout-to-image schemes. We employed the DDIM scheduler [26] with 50 denoising steps in all experiments, performing latent variable optimization only within the first 10 steps, with a maximum of 5 iterations per step. The step size for updating the latent variable  $z_t$  was set to 70, and the loss threshold value for early stopping was set to  $10^{-6}$ . In all experiments, the weight for classifier-free guidance was set to 7.5. In the experiments with MAC,  $\lambda$  and  $\alpha$  for combining boundary-attention loss and regularization loss were both set to 0.25.

## Appendix B: Evaluation Details

We employed the state-of-the-art GroundDINO [15] to detect objects in the synthetic images generated with the input prompts and layout instruction from DrawBench and HRS benchmarks. Specifically, GroundDINO generates multiple predicted bounding boxes corresponding to predicted categories on the synthetic images with threshold value for confidence 0.25. Those boxes are then used to compute various metrics for measuring the spatial controllability and counting accuracy.

### Object Counting

To evaluate the layout-to-image schemes in object counting, we record the number of predicted bounding boxes  $n_{\text{pred}}^{(i)}$  corresponding to the phrase  $\mathbf{p}^{(i)}$ , and compute the correct number of boxes and false number of boxes by

$$n_{\text{cor}}^{(i)} = \min(n_{\text{pred}}^{(i)}, n_{\text{gt}}^{(i)}), \quad (13)$$

$$n_{\text{fal}}^{(i)} = \max(n_{\text{pred}}^{(i)} - n_{\text{gt}}^{(i)}, 0), \quad (14)$$

$$n_{\text{neg}}^{(i)} = \max(n_{\text{gt}}^{(i)} - n_{\text{pred}}^{(i)}, 0), \quad (15)$$

where  $n_{\text{gt}}^{(i)}$  is the ground-truth number. With  $n_{\text{cor}}^{(i)}$  and  $n_{\text{fal}}^{(i)}$ , we can obtain

$$\text{precision} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)} + \sum_{i \in \mathcal{I}} n_{\text{fal}}^{(i)}}, \quad (16)$$

$$\text{recall} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)} + \sum_{i \in \mathcal{I}} n_{\text{neg}}^{(i)}}, \quad (17)$$

where  $\mathcal{I}$  denotes indices of the phrases in the prompt, which is defined in the Section 3.1 in the main paper. With the metrics of precision and recall, we compute  $F1$  score by

$$F1 = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}. \quad (18)$$

### Spatial Accuracy

In the experiments on spatial relationship, phrases in the input prompts are given ground-truth relationship. For instance, in the prompt *a cat is on the left of a dog*, the two phrases *cat* and *dog* have the spatial relationship *on the left of*. By comparing the mean points of the predicted boxes for *cat* and *dog* as

$$(\mathbf{x}_1, \mathbf{y}_1) = \left( \frac{\mathbf{x}_1^{\text{left}} + \mathbf{x}_1^{\text{right}}}{2}, \frac{\mathbf{y}_1^{\text{left}} + \mathbf{y}_1^{\text{right}}}{2} \right), \quad (19)$$

$$(\mathbf{x}_2, \mathbf{y}_2) = \left( \frac{\mathbf{x}_2^{\text{left}} + \mathbf{x}_2^{\text{right}}}{2}, \frac{\mathbf{y}_2^{\text{left}} + \mathbf{y}_2^{\text{right}}}{2} \right), \quad (20)$$

and record the number of correct prediction

$$n_{\text{cor}}^{(i)} = \begin{cases} 1, & \text{if } (\mathbf{x}_1, \mathbf{y}_1) \text{ is on the ground-truth of } (\mathbf{x}_2, \mathbf{y}_2), \\ 0, & \text{if } (\mathbf{x}_1, \mathbf{y}_1) \text{ is not on the ground-truth of } (\mathbf{x}_2, \mathbf{y}_2), \end{cases} \quad (21)$$

and compute the spatial accuracy

$$ACC_{\text{spatial}} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{|\mathcal{I}|}. \quad (22)$$

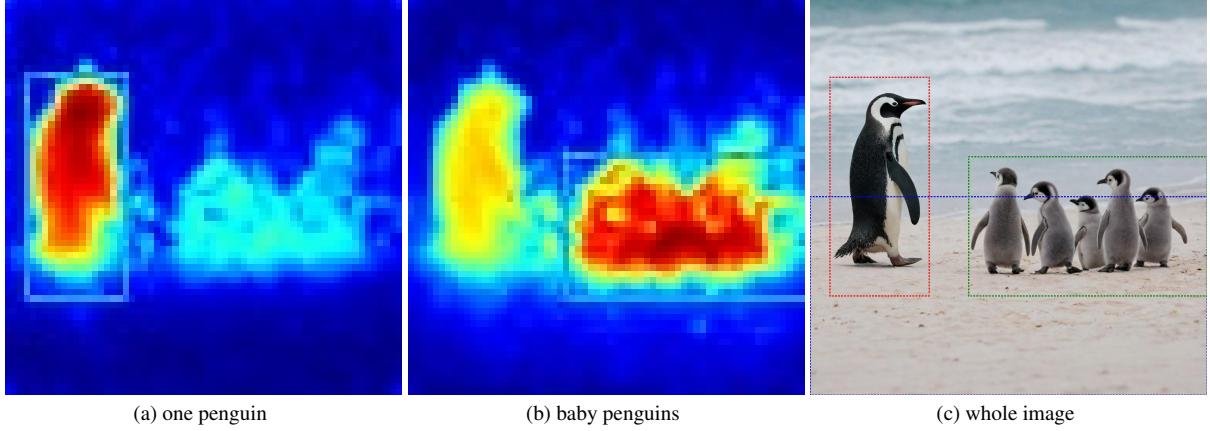


Figure 6. Example of overlapping objects in a box.

## Size Accuracy

In the experiments on size relationship, phrases in the input prompts are given ground-truth relationship. For instance, in the prompt *a cat is smaller than a dog*, the two phrases *cat* and *dog* have the size relationship *smaller*. By comparing the area of the predicted boxes  $\mathcal{B}_1$  and  $\mathcal{B}_2$  for *cat* and *dog* as

$$n_{\text{cor}}^{(i)} = \begin{cases} 1, & \text{if } \text{Area}(\mathcal{B}_1) \text{ is ground-truth than } \text{Area}(\mathcal{B}_2), \\ 0, & \text{if } \text{Area}(\mathcal{B}_1) \text{ is not ground-truth than } \text{Area}(\mathcal{B}_2), \end{cases} \quad (23)$$

and compute the size accuracy

$$ACC_{\text{size}} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{|\mathcal{I}|}. \quad (24)$$

## Color Accuracy

In the experiments on color, phrases in the input prompts are given color instruction. For instance, in the prompt *a white cat and a black dog*, the two phrases *cat* and *dog* have the color *white* and *black*. Similarly, we can compute the color accuracy:

$$n_{\text{cor}}^{(i)} = \begin{cases} 1, & \text{if the predicted color is ground-truth,} \\ 0, & \text{if the predicted color is not ground-truth,} \end{cases} \quad (25)$$

and compute the color accuracy

$$ACC_{\text{color}} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{|\mathcal{I}|}. \quad (26)$$

## Complex Layouts

As shown in Figure 6, our method excels at handling complex and crowded scenes compared to previous approaches, thanks to the boundary-attention constraint. For **overlapping objects**, the user can assign a unique bounding box to each object and use a distinct phrase to describe them. For example: “One penguin is standing on the left of the beach, and five baby penguins are on the right.” The user can assign one box for the phrase **One penguin** while one box for the phrase **baby penguins**, as illustrated on the top right. However, the cross-attention maps for multiple objects within the same box inevitably overlap, leading to incorrect counting. We recognize that these training-free L2I schemes, including ours, may struggle with accurate counting in overlapping boxes, resulting in unsatisfactory generation in crowded scenarios. Addressing this limitation would be interesting for future work.

## Hyper-parameter Sensitivity

The results in Table 4 highlight the roles of the three loss functions: (1) region-attention loss regulates object location, (2) boundary-attention loss improves object counting accuracy, and (3) regularization loss prevents the attention map from vanishing. Therefore, we set  $\lambda = \alpha = 0.5 < 1$  as location regulation is the highest priority. To provide further insights, we will conduct additional experiments with varying  $\lambda$  and  $\alpha$  in the revised version.

## Text-Layout Mismatch

The entire discussion on improving object counting is based on the assumption that the user inputs a consistent layout and text prompts, which is a reasonable assumption, our method will not improve the counting otherwise. Under the same assumption, the existing schemes have shown poor performance in object counting accuracy, as illustrated in Table 1, because it is challenging to accurately determine the count relying solely on the text prompt.