

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

Supplementary Material

Anonymous CVPR submission

Paper ID 11

001 Appendix A: Experimental Details

002 In this section, we report the model architectures and hyper-
003 parameters used in the experiments across all layout-to-
004 image schemes. We employed the DDIM scheduler [2]
005 with 50 denoising steps in all experiments, performing la-
006 tent variable optimization only within the first 10 steps, with
007 a maximum of 5 iterations per step. The step size for updat-
008 ing the latent variable z_t was set to 70, and the loss threshold
009 value for early stopping was set to 10^{-6} . In all experiments,
010 the weight for classifier-free guidance was set to 7.5. In the
011 experiments with MAC, λ and α for combining boundary-
012 attention loss and regularization loss were both set to 0.25.

013 Appendix B: Evaluation Details

014 We employed the state-of-the-art GroundDINO [1] to de-
015 tect objects in the synthetic images generated with the input
016 prompts and layout instruction from DrawBench and HRS
017 benchmarks. Specifically, GroundDINO generates multiple
018 predicted bounding boxes corresponding to predicted cate-
019 gories on the synthetic images with threshold value for confi-
020 dence 0.25. Those boxes are then used to compute various
021 metrics for measuring the spatial controllability and count-
022 ing accuracy.

023 Object Counting

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

$$028 \quad n_{\text{cor}}^{(i)} = \min(n_{\text{pred}}^{(i)}, n_{\text{gt}}^{(i)}), \quad (1)$$

$$029 \quad n_{\text{fal}}^{(i)} = \max(n_{\text{pred}}^{(i)} - n_{\text{gt}}^{(i)}, 0), \quad (2)$$

$$030 \quad n_{\text{neg}}^{(i)} = \max(n_{\text{gt}}^{(i)} - n_{\text{pred}}^{(i)}, 0), \quad (3)$$

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

$$033 \quad \textit{precision} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)} + \sum_{i \in \mathcal{I}} n_{\text{fal}}^{(i)}}, \quad (4) \quad 034$$

$$035 \quad \textit{recall} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)} + \sum_{i \in \mathcal{I}} n_{\text{neg}}^{(i)}}, \quad (5) \quad 036$$

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

$$037 \quad \textit{F1} = \frac{2 * \textit{precision} * \textit{recall}}{\textit{precision} + \textit{recall}}. \quad (6) \quad 038$$

039 Spatial Accuracy

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

$$046 \quad (\mathbf{x}_1, \mathbf{y}_1) = \left(\frac{\mathbf{x}_1^{\text{left}} + \mathbf{x}_1^{\text{right}}}{2}, \frac{\mathbf{y}_1^{\text{left}} + \mathbf{y}_1^{\text{right}}}{2} \right), \quad (7) \quad 047$$

$$048 \quad (\mathbf{x}_2, \mathbf{y}_2) = \left(\frac{\mathbf{x}_2^{\text{left}} + \mathbf{x}_2^{\text{right}}}{2}, \frac{\mathbf{y}_2^{\text{left}} + \mathbf{y}_2^{\text{right}}}{2} \right), \quad (8) \quad 049$$

and record the number of correct prediction

$$050 \quad n_{\text{cor}}^{(i)} = \begin{cases} 1, & \text{if } (\mathbf{x}_1, \mathbf{y}_1) \text{ is on the } \textit{ground-truth} \text{ of } (\mathbf{x}_2, \mathbf{y}_2), \\ 0, & \text{if } (\mathbf{x}_1, \mathbf{y}_1) \text{ is not on the } \textit{ground-truth} \text{ of } (\mathbf{x}_2, \mathbf{y}_2), \end{cases} \quad (9) \quad 051$$

and compute the spatial accuracy

$$052 \quad \textit{ACC}_{\text{spatial}} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{|\mathcal{I}|}. \quad (10) \quad 053$$

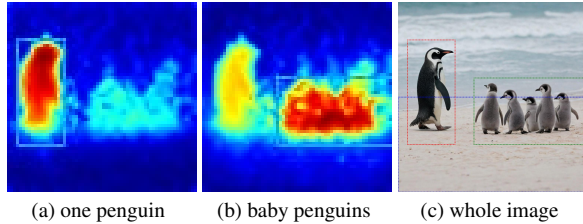


Figure 1. Example of overlapping objects in a box.

056 Size Accuracy

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

$$n_{\text{cor}}^{(i)} = \begin{cases} 1, & \text{if Area}(\mathcal{B}_1) \text{ is } \textit{ground-truth} \text{ than Area}(\mathcal{B}_2), \\ 0, & \text{if Area}(\mathcal{B}_1) \text{ is not } \textit{ground-truth} \text{ than Area}(\mathcal{B}_2), \end{cases} \quad (11)$$

062 and compute the size accuracy

$$064 \quad ACC_{\text{size}} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{|\mathcal{I}|}. \quad (12)$$

065 Color Accuracy

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

$$n_{\text{cor}}^{(i)} = \begin{cases} 1, & \text{if the predicted color is } \textit{ground-truth}, \\ 0, & \text{if the predicted color is not } \textit{ground-truth}, \end{cases} \quad (13)$$

071 and compute the color accuracy

$$073 \quad ACC_{\text{color}} = \frac{\sum_{i \in \mathcal{I}} n_{\text{cor}}^{(i)}}{|\mathcal{I}|}. \quad (14)$$

074 Complex Layouts

075 As shown in Figure 1, our method excels at handling com-
076 plex and crowded scenes compared to previous approaches,
077 thanks to the boundary-attention constraint. For **overlap-**
078 **ping objects**, the user can assign a unique bounding box
079 to each object and use a distinct phrase to describe them.
080 For example: “One penguin is standing on the left of the
081 beach, and five baby penguins are on the right.” The user
082 can assign one box for the phrase **One penguin** while one
083 box for the phrase **baby penguins**, as illustrated on the top

right. However, the cross-attention maps for multiple ob-
jects within the same box inevitably overlap, leading to in-
correct counting. We recognize that these training-free L2I
schemes, including ours, may struggle with accurate count-
ing in overlapping boxes, resulting in unsatisfactory genera-
tion 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 loca-
tion, (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 fur-
ther insights, we will conduct additional experiments with
varying λ and α 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.

References

- [1] 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. 1
- [2] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 1