ReGround: Improving Textual and Spatial Grounding at No Cost – Supplementary Material

Phillip Y. Lee[®] Minhyuk Sung[®]

KAIST {phillip0701,mhsung}@kaist.ac.kr

In this supplementary material, we provide additional details on the evaluation setup (Sec. S.1) and more quantitative comparisons of our REGROUND and GLIGEN [3] (Sec. S.2). Moreover, we showcase the effect of REGROUND as a backbone of zero-shot layout-guided image generation methods (Sec. S.3) and finally provide extensive qualitative comparisons of Stable Diffusion (SD) [8], GLIGEN [3], and our REGROUND (Sec. S.4).

S.1 Details on Evaluation Setup

This section provides further descriptions on the evaluation datasets (Sec. 6.1) and the user study setup (Sec. 6.2) in the main paper

MS-COCO. The validation set of the MS-COCO-2017 dataset [5] consists of 5,000 image-annotation pairs. Since GLIGEN [3] is trained to handle a maximum of 30 bounding boxes per image, we excluded pairs with more than 30 bounding boxes or no bounding boxes, resulting in a total of 4,952 images. For the validation set of the MS-COCO-2014 dataset [5], we randomly sampled 5,000 pairs for evaluation.

NSR-1K-GPT. Numerical and Spatial Reasoning (NSR-1K) [1] is a collection of layout-caption pairs designed to assess the numerical and spatial reasoning capabilities of image generation methods. The object labels and bounding boxes are from MS-COCO [5], while the captions are newly annotated based on the spatial relationships and numerical properties of objects. NSR-1K consists of two subsets: *Counting* and *Spatial*. We randomly sampled 1,000 pairs from the Counting set and used all 1,021 pairs from the Spatial set.

User Study. We conducted the user study through Amazon Mechanical Turk using the template displayed in Fig. S1. Based on the text prompt and bounding boxes generated from GPT-4 [6], images were generated by both GLIGEN [3] and our REGROUND. Since REGROUND aims to resolve the failure cases of GLIGEN, we re-generated both images when the differences between them were minimal (*i.e.*, if the LPIPS value [11] was less than 0.3), resulting in an average of 2.4 iterations per image. Each participant answered 20 questions and 5 vigilance tests.



Fig. S1: User study template. In the above example, the text prompt "A photo of a bicycle and a bench" was displayed to the respondents.

S.2 Additional Quantitative Comparisons

In addition to Sec. 6.3 of the main paper, this section provides quantitative comparisons between GLIGEN [3] and our REGROUND on the Spatial subset of NSR-1K-GPT, and with a different version of Stable Diffusion [8] as the base image diffusion model.

Comparison on NSR-1K-GPT-Spatial. Fig. S2-(a) shows the CLIP score [2] and YOLO score [9] measured on the *Spatial* subset of NSR-1K-GPT. The minimum CLIP score of our REGROUND (33.89 at $\gamma = 1.0$) is already higher than GLIGEN's maximum CLIP score (33.88 at $\gamma = 0.1$), indicating that REGROUND obtains a significant enhancement in textual grounding while preserving the spatial grounding.



Fig. S2: Quantitative comparisons (a) on the *Spatial* subset of NSR-1K-GPT and (b) using SDv2.1 as the base image diffusion model. Consistent with the findings from Fig. 6 of the main paper, our REGROUND demonstrates improved performance in textual and spatial groundings, as seen by the higher CLIP score [2] for the same range of YOLO score [9].

Results with SDv2.1 as Base Diffusion Model. In Sec. 6 of the main paper, we conducted experiments using the GLIGEN [3] checkpoint based on Stable Diffusion v1.4 (SDv1.4). Additionally, we provide quantitative comparisons with an unofficial GLIGEN checkpoint [4] that was trained with SDv2.1 as the base image diffusion model. The results, presented in Fig. S2-(b), clearly demonstrate the significant outperformance of our REGROUND over GLIGEN.

S.3 More Results with REGROUND as Backbone

In addition to Sec. 6.4 of the main paper, we provide qualitative comparisons of different layout-guided generation methods using GLIGEN [3] and our RE-GROUND as backbones, respectively (Fig. S3, S4). The results on BoxDiff [10] and Attention Refocusing [7] illustrate that our network rewiring substantially improves the performance of layout-guided generation methods built upon the GLIGEN framework.



"A cat in a wooden room wearing a birthday hat."

Fig. S3: Comparisons of GLIGEN [3] and our REGROUND as a backbone for BoxDiff [10] and Attention Refocusing (Attn-Refocus) [7].



"A <u>black and white photo</u> of an apple and a cup on a table."

Fig. S4: More comparisons on BoxDiff [10] and Attention Refocusing (Attn-Refocus) [7].

S.4Additional Qualitative Comparisons

In this section, we provide extensive qualitative comparisons of Stable Diffusion (SD) [8], GLIGEN [3], and our REGROUND on layout-guided image generation. Note that $\gamma \in [0,1]$ denotes the fraction of the initial denoising steps during which gated self-attention is activated, as discussed in Sec. 5.1 of the main paper.

In each row, the input layout is presented in the first column, with the input text prompt displayed below the images. The phrase <u>underlined</u> in each prompt highlights the entity subject to description omission, as mentioned in Sec. 4.2 of the main paper. Furthermore, black arrows are used to denote bounding boxes that some methods fail to represent accurately, whereas other methods succeed in doing so precisely. Red arrows signify a failure in either spatial or textual grounding, while green arrows indicate successful grounding of a specific entity.



"A person covered with an orange sleeping bag [...] sleeps on a park bench."



"A piece of cake that is sitting on foil.



"This is a picture of <u>an interesting mural</u> on a truck.





"A person standing next to a green and black train.



"A tan dog eating food scraps from a plate.



"A room with a fireplace and television inside of it."



"[...] a toilet with a gleaming gold lid and ornate rim around it stands as a symbol of opulence...



"A wall shelf with a single, delicate vase and multiple vintage picture frames, [...].



"<u>A mirror</u> is hanging next to a vase of flowers."



"A dish on the counter covered with aluminum foil."



References

- Feng, W., Zhu, W., Fu, T.j., Jampani, V., Akula, A., He, X., Basu, S., Wang, X.E., Wang, W.Y.: LayoutGPT: Compositional visual planning and generation with large language models. In: NeurIPS (2023) 1
- 2. Hessel, J., Holtzman, A., Forbes, M., Bras, R.L., Choi, Y.: CLIPScore: a reference-free evaluation metric for image captioning. In: EMNLP (2021) 2
- 3. Li, Y., Liu, H., Wu, Q., Mu, F., Yang, J., Gao, J., Li, C., Lee, Y.J.: GLIGEN: Open-set grounded text-to-image generation. In: CVPR (2023) 1, 2, 3, 5
- Lian, L., Li, B., Yala, A., Darrell, T.: Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models. arXiv preprint arXiv:2305.13655 (2023) 3
- Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014) 1
 Oper Al: Chatrant https://chat.approx.org/1
- 6. OpenAI: Chatgpt, https://chat.openai.com/ 1
- 7. Phung, Q., Ge, S., Huang, J.B.: Grounded text-to-image synthesis with attention refocusing. arXiv preprint arXiv:2306.05427 (2023) 3, 4
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR (2022) 1, 2, 5
- 9. Wang, C.Y., Bochkovskiy, A., Liao, H.Y.M.: YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. In: CVPR (2023) 2
- Xie, J., Li, Y., Huang, Y., Liu, H., Zhang, W., Zheng, Y., Shou, M.Z.: Boxdiff: Textto-image synthesis with training-free box-constrained diffusion. In: ICCV (2023) 3, 4
- 11. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: CVPR (2018) 1