# Appendix

In this supplemental file, we provide the following materials:

- Sec. A.1: more illustrations of noise prediction error  $\epsilon_{FT}(t)$  by different diffusion models  $\epsilon(t)$  (referring to Sec. 3.1 and Fig. 1 in the main paper);
- Sec. A.2: more 2D toy experiments of different methods (referring to Sec. 3.2 and Fig. 2 in the main paper);
- Sec. A.3: more details of 3D generator architectures (referring to Sec. 3.2 and Fig. 3 in the main paper);
- Sec. A.4: more corpus details (referring to Sec. 4.1 in the main paper);
- Sec. A.5: more implementation details (referring to Sec. 4.2 in the main paper);
- Sec. A.6: more results with MVDream as the 2D diffusion prior. Including the prompt-specifc generation and prompt-amortized training (referring to Sec. 4.2 and Sec. 4.3 in the main paper);
- Sec. A.7: the discussion and comparison with datadriven methods (referring to Sec. 4.2 in the main paper);

## A.1 More Illustrations of Noise Prediction Error

In this section, we provide more illustrations of the noise prediction error by various pre-trained diffusion models, including the 2D  $\epsilon$ -prediction model [2, 29] and the v-prediction model [1, 30], and the 3D diffusion model [6]. We plot the the noise prediction error against timesteps in Fig. 1. For each text prompt displayed at the top of the sub-figures, we use it as the condition to generate 16 samples. We then introduce a single instance of Gaussian noise to each sample and execute one diffusion step at 100 different timesteps. The DDPM [9] is used as the noise scheduler, as done in VSD [38]. The average noise reconstruction error is then calculated over the timesteps and the 16 data samples.

2D  $\epsilon$ -prediction diffusion model. The  $\epsilon$ -prediction model is widely adopted in the field of text-to-3D synthesis [17,28,33,38,41]. In our tests, we employ the commonly used SD-v2.1-base model [2]. The noise prediction error curves for four prompts sourced from Magic3D [18] are presented in Fig. 1(a), from which we see a clear decrease of noise prediction error with the timestep going from  $T_{\rm min}$  to  $T_{\rm max}$ .

**2D** *v*-prediction diffusion model. The *v*-prediction model, introduced by Salimans *et al.* [30], accelerates the generation process by predicting velocity rather than noise. We test this model using the well-known SD-v2.1 [1] with 4 prompts sourced from Magic3D [18]. To calculate the noise prediction error, we convert the velocity predictions into noise predictions [30]. As depicted in Fig. 1(b), the *v*-prediction model also exhibits reduced prediction errors as the timestep goes from  $T_{\rm min}$  to  $T_{\rm max}$ .

**3D** diffusion model. Apart from the above 2D diffusion models, we also conduct experiments on a 3D diffusion model DiffTF [6], which is a 3D generator trained on 3D object datasets [40]. It is configured with  $\epsilon$ -prediction and



**Fig. 1:** The behavior of noise prediction error of different diffusion models, including (a) 2D  $\epsilon$ -prediction [2] diffusion model, (b) 2D *v*-prediction [1] diffusion model, and (c) 3D diffusion model. Zoom in for a better view.

performs the diffusion process on tri-plane [7]. As shown in Fig. 1(c), its noise prediction error e(t) also reduces as timestep t increases, which is similar to 2D diffusion models. In particular, e(t) drops rapidly before t = 200. This is mainly caused by the much smaller scale (e.g., 6k 3D objects) of the 3D dataset [8] compared with the 2D datasets [32] (e.g., 2B text-image pairs). Therefore, the network tends to overfit the 3D data with smaller prediction error.

## A.2 More 2D Toy Experiments

To further validate the effectiveness of the introduced timestep interval  $\Delta t$  in our ASD, we provide more 2D toy experiments in Fig. 2, covering a wild range of subjects, *i.e.*, plants, objects, animals, and scenes.

From Fig. 2, we can see that SDS [26] and CSD [47] do not perform very well. SDS generates high-saturation results because of the large CFG [10], while CSD shows noisy and blurred patterns so that the subjects are difficult to identify. VSD generates good quality results by fine-tuning the 2D diffusion model. However, as we discussed in the main paper, it hurts the 2D diffusion model's comprehension capability to numerous text prompts, leading to mode collapse



Fig. 2: 2D toy experiments by SDS [26], CSD [10], VSD [38] and our ASD with different settings of  $\Delta t$ .

when the size of text prompts is extended. Without changing the diffusion prior, our proposed ASD can achieve the same high quality results as VSD.

We also ablate the setting of  $\Delta t$  in this experiment. We see that if we set  $\Delta t = 0$ , it leads to a noisy pattern similar to CSD. By setting it as a fixed interval, *e.g.*,  $\Delta t = \eta T_{\text{max}}$ , it would result in poor texture or geometry, such as the panda in Fig. 2. By setting  $\Delta t$  relevant to t as  $\Delta t = \eta (t - T_{\min})$ , the results can be much improved. Finally, the results are further enhanced by randomly sampling  $\Delta t$  via  $\Delta t \sim \mathcal{U}[0, \eta (t - T_{\min})]$ . The detailed explanations can be found in Sec. 3.2 of the main paper.

## A.3 More 3D Generator Architecture Details

Hyper-iNGP. We replicate the hypernetwork design from ATT3D [21], integrating it with iNGP [25] to achieve prompt-amortized text-to-3D synthesis. As illustrated in Fig. 3, the hypernetwork projects text prompt embeddings into the



**Fig. 3:** The network architecture and rendering scheme of Hyper-iNGP(left), 3DConvnet(middle) and Triplane-Transformer (right)

weights of linear layers. The HashGrid representation [25] encodes sample points independently, which are then transformed by the hypernetwork-parameterized linear layers into prompt-specific color c and density  $\sigma$ . Following ATT3D [21], another hypernetwork is implemented to create a prompt-specific background. The ray direction is encoded into a separate HashGrid, which is then projected to the background color  $c_{bg}$ , facilitating the creation of high-resolution backgrounds. The spectral normalization [24] can be optionally turned on to stabilize the training with SDS [26].

**3DConv-net**. As illustrated in Fig. 3, our 3DConv-net mirrors the Style-GAN2 model [14], using modulated convolutions to upscale features directed by the latent code  $\mathbf{w}$ , which is conditioned on Gaussian noise  $\mathbf{z} \sim \mathcal{N}(0, 1)$  and the text prompt embedding as in text-driven 2D GANs [31]. Transitioning from 2D to 3D, we substitute StyleGAN2's components with their 3D alternatives, modulated by  $\mathbf{w}$ . The network up-samples a 4<sup>3</sup> dimensional voxel to 128<sup>3</sup> dimension. For quicker convergence, we add 3D bias within blocks for processing voxels with the dimension from 8<sup>3</sup> to 64<sup>3</sup>. Rendering is accomplished by interpolating voxel features to determine the color and density of each point along the rays. A background module is incorporated as well.

**Triple-Transformer**. Recently, the Transformer [37] architecture has gained popularity in 3D generation tasks for its scalability, especially in data-driven methods [12, 15, 20, 36, 39, 43–45, 49]. However, it has not been applied in recent score-distillation-based methods yet [16, 27, 42]. In this paper, we conduct experiments to explore the performance of Transformer architecture in score-distillation-based text-to-3D generation. As shown in Fig. 3, we employ 12 Transformer layers, each comprising self-attention, cross-attention, and feed-forward networks. The text prompt is first processed by the CLIP text encoder and then fed into the cross-attention to set the condition. The query embeddings



**Fig. 4:** Qualitative comparison between SDS\* and ASD on prompt-specific text-to-3D generation, with iNGP as 3D representation and MVDream as 2D diffusion prior.

are passed through these layers, and then reshaped and up-sampled to form a triplane, which is an efficient 3D representation [7].

**Rendering**. For prompt-specific optimization, we use the volume rendering in NeRF [38] and keep the configuration in prior arts [38]. For prompt-amortized training, we implement VolSDF [46], which uses 64 sample points for coarse sampling and 256 sample points for fine sampling [23]. We found that keeping the mean absolute deviation fixed to be 30 can achieve good results. We render  $64 \times 64$  resolution for 3DConv-net and  $256 \times 256$  for Hyper-iNGP in the whole training period.

#### A.4 More Details about Corpus

In this work, we utilize five corpora to assess our ASD for prompt-based text-to-3D generation. Apart from MG15 [18], DF415 [26], AT2520 [21] and DL17k [16], we also provide the CP100k corpus. CP100k consists of 100k corpus for training and 1k corpus for test, which are sampled from Cap3D [22].

## A.5 More Implementation Details

**Prompt-specific Text-to-3D**. Our code is based on the open-source Text-to-3D codebase [3]. We follow the configuration in ProlificDreamer [4] in specifying the parameters, including the training iterations, optimizer, batch-size and learning rate. All experiments are conducted on one Nvidia V100 GPU.

**Prompt-amortized Text-to-3D**. The experiments for prompt-amortized text-to-3D are conducted on 8 Nvidia A6000 GPUs, with a per-GPU batch size of 1. Training on MG15, DF417, AT2520, DL17k and CP100k requires 50k, 100k, 50k, 200k and 300k iterations, respectively.

**2D Diffusion Guidance**. For 2D experiments, utilizing the diffusion model [2] with T = 1000 timesteps, we adhere to the existing protocol [4] by setting  $T_{\min} = 20$  and  $T_{\max} = 980$ . In the 3D experiments, we adopt the approaches in [38] and [33], where  $T_{\max}$  is progressively reduced from 980 to 500 to enhance the quality of generation outputs. We start with a higher  $T_{\min}$  and decrease it



"A petite woman with glasses is typing on her laptop in a bustling café"  ${\rm in}\,{\rm DL17k}$ 

Fig. 5: Qualitative comparison among SDS\* [33] and our ASD on DL17k corpus with Triplane-Transformer as 3D generator and MVDream as 2D diffusion prior.

|         |  |              | $\operatorname{Sim}\uparrow$ | $R@1\uparrow$ |
|---------|--|--------------|------------------------------|---------------|
| $SDS^*$ |  |              | 0.200                        | 0.159         |
|         | $\Delta t = \eta (t - T_{\min})$                   | $\eta = 0.1$ | 0.205                        | 0.231         |
| ASD     | $\Delta t \sim 1/[0 m(t - T)]$                     | $\eta = 0$   | 0.213                        | 0.293         |
|         | $\Delta t \sim \mathcal{U}[0, \eta(t - I_{\min})]$ | $\eta = 0.1$ | 0.219                        | 0.294         |

**Table 1:** Comparison with SDS\* and ablation study on ASD using MVDream as the2D diffusion model.

linearly from 500 to 20, which helps to mitigate the Janus issue, as adopted in [5]. Additionally, when Stable Diffusion is used as the 2D diffusion model, we employ the Perp-neg strategy [5] to further address the Janus problem.

# A.6 Results with MVDream

As a score distillation method, ASD is open to the choice of 2D diffusion models. In this section, we evaluate ASD's compatibility with another representative 2D diffusion model, MVDream [33]. To conduct score distillation, MVDream takes four views as input for rendering, and explicitly uses the camera poses as prompts. We conduct comparison and ablation study in prompt-specific optimization with iNGP as the 3D representation, as well as prompt-amortized text-to-3D with Triplane-Transformer as the 3D generator.

**Results with iNGP as 3D Representation**. MVDream officially implements a modified SDS method by incorporating the CFG re-scale technique [19] to alleviate large gradient norms caused by SDS. We refer to this modified SDS as SDS\*. We qualitatively compare the performance of SDS\* and ASD on prompt-specific text-to-3D. The results are shown in Fig. 4. It can be seen that SDS\* produces abnormal geometry with solid matter covering most of the 3D space,



"A dog made out of salad" "A baby dragon" Fig. 6: The visual comparison with data-driven methods LGM [34] and Shape-E [13].

and it generates grayish textures. In contrast, ASD generates more natural geometry and textures.

Results with Triplane-Transformer as 3D Generator. We then employ MVDream for prompt-amortized text-to-3D by using Triplane-Transformer as the 3D generator. In addition to the comparison with SDS\*, we ablate ASD without timestep shift to further solidify our proposed asynchronous timesteps. The experiments are conducted on DL17k corpus. As shown in Fig. 5, SDS\* tends to produce small geometries. By using ASD with a deterministic timestep shift, *i.e.*  $\Delta t = \eta (t - T_{\min})$ , the results are improved yet still unsatisfactory. Without any timestep shift in ASD, *i.e.*,  $\eta = 0$ , the 3D results have some floating patterns. This happens because without a timestep shift, the model fails to align the distribution of rendered images with the prior distribution of pre-trained diffusion model. By using a random timestep shift  $\Delta t \sim \mathcal{U}[0, \eta (t - T_{\min})]$  and the magnitude of  $\eta = 0.1$  in ASD, the results are significantly improved, which is also reflected in the metrics shown in Tab. 1.

## A.7 Discussions with Data-Driven Methods

Our proposed method differs from existing data-driven methods [11,13,34,35,48] in that we do not require any 3D dataset to train the 3D generator. If the test text prompts fall into the training distribution, these supervised data-driven methods may generate better quality outputs than our unsupervised method. However, by leveraging the strong prior information in pre-trained 2D diffusion models, our method has better generalization capability to the test prompts. By using our 3DConv-net trained on DF415 corpus as an example, we compare our results with open-sourced data-driven 3D generators LGM [34] and Shape-E [13]. Fig. 6 shows the qualitative comparison on some text prompt inputs, which are are out of the training distribution. We can see that LGM and Shape-E output poor results. In contrast, ASD can still work well by exploiting the powerful diffusion priors in pre-trained 2D models.

## References

- 1. Stable-diffusion-v2.1. https://huggingface.co/stabilityai/stable-diffusion-2-1
- Stable-diffusion-v2.1-base. https://huggingface.co/stabilityai/stablediffusion-2-1-base
- 3. Threestudio: a unified framework for 3d content creation from text prompts. https: //github.com/threestudio-project/threestudio
- Unofficial implementation of 2d prolificdreamer. https://github.com/yuanzhizhu/prolific\_dreamer2d
- Armandpour, M., Zheng, H., Sadeghian, A., Sadeghian, A., Zhou, M.: Re-imagine the negative prompt algorithm: Transform 2d diffusion into 3d, alleviate janus problem and beyond. arXiv preprint arXiv:2304.04968 (2023)
- Cao, Z., Hong, F., Wu, T., Pan, L., Liu, Z.: Large-vocabulary 3d diffusion model with transformer. arXiv preprint arXiv:2309.07920 (2023)
- Chan, E.R., Lin, C.Z., Chan, M.A., Nagano, K., Pan, B., De Mello, S., Gallo, O., Guibas, L.J., Tremblay, J., Khamis, S., et al.: Efficient geometry-aware 3d generative adversarial networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16123–16133 (2022)
- Deitke, M., Liu, R., Wallingford, M., Ngo, H., Michel, O., Kusupati, A., Fan, A., Laforte, C., Voleti, V., Gadre, S.Y., et al.: Objaverse-xl: A universe of 10m+ 3d objects. arXiv preprint arXiv:2307.05663 (2023)
- Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in neural information processing systems 33, 6840–6851 (2020)
- Ho, J., Salimans, T.: Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598 (2022)
- Hong, F., Tang, J., Cao, Z., Shi, M., Wu, T., Chen, Z., Wang, T., Pan, L., Lin, D., Liu, Z.: 3dtopia: Large text-to-3d generation model with hybrid diffusion priors. arXiv preprint arXiv:2403.02234 (2024)
- Hong, Y., Zhang, K., Gu, J., Bi, S., Zhou, Y., Liu, D., Liu, F., Sunkavalli, K., Bui, T., Tan, H.: Lrm: Large reconstruction model for single image to 3d. arXiv preprint arXiv:2311.04400 (2023)
- Jun, H., Nichol, A.: Shap-e: Generating conditional 3d implicit functions. arXiv preprint arXiv:2305.02463 (2023)
- 14. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8110–8119 (2020)
- Li, M., Long, X., Liang, Y., Li, W., Liu, Y., Li, P., Chi, X., Qi, X., Xue, W., Luo, W., et al.: M-lrm: Multi-view large reconstruction model. arXiv preprint arXiv:2406.07648 (2024)
- Li, M., Zhou, P., Liu, J.W., Keppo, J., Lin, M., Yan, S., Xu, X.: Instant3d: Instant text-to-3d generation. arXiv preprint arXiv:2311.08403 (2023)
- Liang, Y., Yang, X., Lin, J., Li, H., Xu, X., Chen, Y.: Luciddreamer: Towards high-fidelity text-to-3d generation via interval score matching. arXiv preprint arXiv:2311.11284 (2023)
- Lin, C.H., Gao, J., Tang, L., Takikawa, T., Zeng, X., Huang, X., Kreis, K., Fidler, S., Liu, M.Y., Lin, T.Y.: Magic3d: High-resolution text-to-3d content creation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 300–309 (2023)
- Lin, S., Liu, B., Li, J., Yang, X.: Common diffusion noise schedules and sample steps are flawed. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 5404–5411 (2024)

8

- Liu, Z., Li, Y., Lin, Y., Yu, X., Peng, S., Cao, Y.P., Qi, X., Huang, X., Liang, D., Ouyang, W.: Unidream: Unifying diffusion priors for relightable text-to-3d generation. arXiv preprint arXiv:2312.08754 (2023)
- Lorraine, J., Xie, K., Zeng, X., Lin, C.H., Takikawa, T., Sharp, N., Lin, T.Y., Liu, M.Y., Fidler, S., Lucas, J.: Att3d: Amortized text-to-3d object synthesis. arXiv preprint arXiv:2306.07349 (2023)
- Luo, T., Rockwell, C., Lee, H., Johnson, J.: Scalable 3d captioning with pretrained models. Advances in Neural Information Processing Systems 36 (2024)
- Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM 65(1), 99–106 (2021)
- Miyato, T., Kataoka, T., Koyama, M., Yoshida, Y.: Spectral normalization for generative adversarial networks. arXiv preprint arXiv:1802.05957 (2018)
- Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. ACM Transactions on Graphics (ToG) 41(4), 1– 15 (2022)
- Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3d using 2d diffusion. arXiv preprint arXiv:2209.14988 (2022)
- 27. Qian, G., Cao, J., Siarohin, A., Kant, Y., Wang, C., Vasilkovsky, M., Lee, H.Y., Fang, Y., Skorokhodov, I., Zhuang, P., et al.: Atom: Amortized text-to-mesh using 2d diffusion. arXiv preprint arXiv:2402.00867 (2024)
- Qiu, L., Chen, G., Gu, X., Zuo, Q., Xu, M., Wu, Y., Yuan, W., Dong, Z., Bo, L., Han, X.: Richdreamer: A generalizable normal-depth diffusion model for detail richness in text-to-3d. arXiv preprint arXiv:2311.16918 (2023)
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022)
- Salimans, T., Ho, J.: Progressive distillation for fast sampling of diffusion models. arXiv preprint arXiv:2202.00512 (2022)
- Sauer, A., Karras, T., Laine, S., Geiger, A., Aila, T.: Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis. arXiv preprint arXiv:2301.09515 (2023)
- 32. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5b: An open largescale dataset for training next generation image-text models. Advances in Neural Information Processing Systems 35, 25278–25294 (2022)
- Shi, Y., Wang, P., Ye, J., Long, M., Li, K., Yang, X.: Mvdream: Multi-view diffusion for 3d generation. arXiv preprint arXiv:2308.16512 (2023)
- 34. Tang, J., Chen, Z., Chen, X., Wang, T., Zeng, G., Liu, Z.: Lgm: Large multiview gaussian model for high-resolution 3d content creation. arXiv preprint arXiv:2402.05054 (2024)
- Tang, Z., Gu, S., Wang, C., Zhang, T., Bao, J., Chen, D., Guo, B.: Volumediffusion: Flexible text-to-3d generation with efficient volumetric encoder. arXiv preprint arXiv:2312.11459 (2023)
- Tochilkin, D., Pankratz, D., Liu, Z., Huang, Z., Letts, A., Li, Y., Liang, D., Laforte, C., Jampani, V., Cao, Y.P.: Triposr: Fast 3d object reconstruction from a single image. arXiv preprint arXiv:2403.02151 (2024)
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)

- Wang, Z., Lu, C., Wang, Y., Bao, F., Li, C., Su, H., Zhu, J.: Prolificdreamer: Highfidelity and diverse text-to-3d generation with variational score distillation. arXiv preprint arXiv:2305.16213 (2023)
- Wei, X., Zhang, K., Bi, S., Tan, H., Luan, F., Deschaintre, V., Sunkavalli, K., Su, H., Xu, Z.: Meshlrm: Large reconstruction model for high-quality mesh. arXiv preprint arXiv:2404.12385 (2024)
- 40. Wu, T., Zhang, J., Fu, X., Wang, Y., Ren, J., Pan, L., Wu, W., Yang, L., Wang, J., Qian, C., et al.: Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 803–814 (2023)
- Wu, Z., Zhou, P., Yi, X., Yuan, X., Zhang, H.: Consistent3d: Towards consistent high-fidelity text-to-3d generation with deterministic sampling prior. arXiv preprint arXiv:2401.09050 (2024)
- 42. Xie, K., Lorraine, J., Cao, T., Gao, J., Lucas, J., Torralba, A., Fidler, S., Zeng, X.: Latte3d: Large-scale amortized text-to-enhanced3d synthesis. arXiv preprint arXiv:2403.15385 (2024)
- 43. Xu, J., Cheng, W., Gao, Y., Wang, X., Gao, S., Shan, Y.: Instantmesh: Efficient 3d mesh generation from a single image with sparse-view large reconstruction models. arXiv preprint arXiv:2404.07191 (2024)
- 44. Xu, Y., Shi, Z., Yifan, W., Chen, H., Yang, C., Peng, S., Shen, Y., Wetzstein, G.: Grm: Large gaussian reconstruction model for efficient 3d reconstruction and generation. arXiv preprint arXiv:2403.14621 (2024)
- 45. Xu, Y., Tan, H., Luan, F., Bi, S., Wang, P., Li, J., Shi, Z., Sunkavalli, K., Wetzstein, G., Xu, Z., et al.: Dmv3d: Denoising multi-view diffusion using 3d large reconstruction model. arXiv preprint arXiv:2311.09217 (2023)
- Yariv, L., Gu, J., Kasten, Y., Lipman, Y.: Volume rendering of neural implicit surfaces. Advances in Neural Information Processing Systems 34, 4805–4815 (2021)
- Yu, X., Guo, Y.C., Li, Y., Liang, D., Zhang, S.H., Qi, X.: Text-to-3d with classifier score distillation. arXiv preprint arXiv:2310.19415 (2023)
- Zhang, B., Cheng, Y., Yang, J., Wang, C., Zhao, F., Tang, Y., Chen, D., Guo, B.: Gaussiancube: Structuring gaussian splatting using optimal transport for 3d generative modeling. arXiv preprint arXiv:2403.19655 (2024)
- 49. Zou, Z.X., Yu, Z., Guo, Y.C., Li, Y., Liang, D., Cao, Y.P., Zhang, S.H.: Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. arXiv preprint arXiv:2312.09147 (2023)

10