# Learning Effective NeRFs and SDFs Representations with 3D Generative Adversarial Networks for 3D Object Generation: Technical Report for ICCV 2023 OmniObject3D Challenge

Zheyuan Yang, Yibo Liu, Guile Wu, Tongtong Cao, Yuan Ren, Yang Liu, Bingbing Liu
Huawei Noah's Ark Lab

{andrew.yang1, yibo.liu3, guile.wu, caotongtong, yuan.ren3, yang.liu9,
liu.bingbing}@huawei.com

## Abstract

*In this technical report, we present a solution for 3D object generation of ICCV 2023 OmniObject3D Challenge. In recent years, 3D object generation has made great process and achieved promising results, but it remains a challenging task due to the difficulty of generating complex, textured and high-fidelity results. To resolve this problem, we study learning effective NeRFs and SDFs representations with 3D Generative Adversarial Networks (GANs) for 3D object generation. Specifically, inspired by recent works, we use the efficient geometry-aware 3D GANs as the backbone incorporating with label embedding and color mapping, which enables to train the model on different taxonomies simultaneously. Then, through a decoder, we aggregate the resulting features to generate Neural Radiance Fields (NeRFs) based representations for rendering high-fidelity synthetic images. Meanwhile, we optimize Signed Distance Functions (SDFs) to effectively represent objects with 3D meshes. Besides, we observe that this model can be effectively trained with only a few images of each object from a variety of classes, instead of using a great number of images per object or training one model per class. With this pipeline, we can optimize an effective model for 3D object generation. This solution is one of the final top-3-place solutions in the ICCV 2023 OmniObject3D Challenge.*

## 1. Introduction

3D object generation aims at generating meaningful 3D surfaces and synthesis images of 3D objects given random inputs [23, 7]. Inspired by the great success of 2D object generation [9, 4], there have been many studies [12, 5] extending 2D methods to 3D generation. Recent 3D generation approaches focus more on effective 3D representation learning [18, 15] and advanced generation strategy [10, 2, 21]. Despite the great progress and promising

results in recent years, 3D object generation is still a challenging task due to the difficulty of generating complex, textured and high-fidelity results.

To devise a solution for 3D object generation in the ICCV 2023 OmniObject3D Challenge, we simultaneously consider effective 3D representation learning and advanced generation strategy. In terms of 3D representation learning, textured mesh [3] and Neural Radiance Fields (NeRFs) [14] embrace the great potential to effectively encode 3D object features. In particular, textured mesh representations can support individual texture map generation, which allows us to easily replace surface textures, while NeRFs-based representations excel at rendering high-fidelity images. Furthermore, when considering effective shape representation, there are a variety of options, ranging from explicit representations like point-based [20], mesh-based [1] and voxel-based [22] approaches to implicit ones such as the Signed Distance Functions (SDFs) [17, 11] and occupancy function [13]. In the light of these cutting-edge techniques, we mainly focus on exploring effective NeRFs representations for 3D object generation while still optimizing SDFs representations to effectively represent objects with 3D meshes.

Moving on to advanced generation strategy, there are usually different options, such as latent generation [21] and direct generation [16]. In particular, latent generation methods aim to generate the representation of the desired output in the latent space and decode the latent code back to the 3D space, while direct generation focus on directly generating the desired 3D representation without utilizing the latent space. Besides, from the perspective of generation process, generation strategies can be categorized into Generative Adversarial Networks (GANs) [10, 2], Variational AutoEncoder (VAE) [19], flow-based [6], Diffusion Models (DMs) [21], *etc*. In our solution, we employ the efficient geometry-aware 3D GANs [2] as the backbone for 3D object generation.

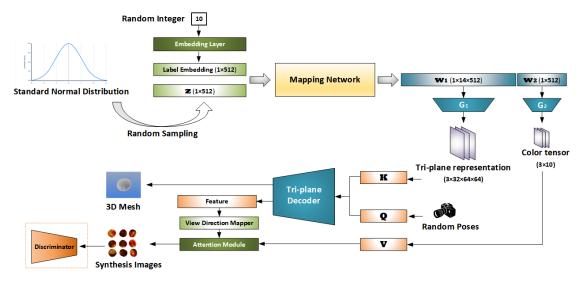The overall framework of our solution is depicted in

Figure 1. An overview of the framework of our solution for 3D object generation.

Fig. 1. Specifically, we aim to learn effective NeRFs and SDFs representations with 3D Generative Adversarial Networks (GANs) for 3D object generation. Inspired by the recent success of efficient geometry-aware 3D GANs, *i.e.*, EG3D [2, 18], we employ it as the backbone and incorporate label embedding and color mapping [18] to enable model training on different taxonomies simultaneously. Next, we generate NeRFs-based representations for rendering high-fidelity synthetic images and SDFs to effectively represent objects with 3D meshes through a decoder. In addition, we empirically find that we can effectively train this model with only a few images per object from various classes, rather than training with a great number of images per object or optimizing one model per class. With this pipeline, we can optimize an effective model for 3D object generation. This solution is one of the final top-3-place solutions in the ICCV 2023 OmniObject3D Challenge.

## 2. Technical Solution

As aforementioned, to realize 3D object generation for ICCV 2023 OmniObject3D Challenge, we explore 3D GANs to learn effective NeRFs representations for synthetic images and SDFs representations for 3D shapes.

**Motivation of Using GANs instead of DMs.** Although diffusion models have recently gained an increasing popularity in 3D object generation, most existing works [16, 21] have only been verified to work under specific category conditions (e.g., chairs [16] or avatars [21]). Besides, since the OmniObject3D dataset [23] is a collection of objects from hundreds of taxonomies, following the pipeline of existing DMs would necessitate the use of hundreds of pre-trained models, which is impractical. Hence, we employ a method

that enables simultaneous training on all the taxonomies in the OmniObject3D dataset. Moreover, given the success of StyleGAN2 [10] in representing NeRFs in latent space, we have chosen latent generation for our pipeline.

**Solution Pipeline.** The framework of our solution is depicted in Fig. 1. Overall, our solution adopts a GAN-based generation process. Specifically, with EG3D [2] as the backbone, firstly, we generate a one-dimensional latent code, $\mathbf{z}$ $(1 \times 512)$, which is obtained by randomly sampling 512 values from a standard normal distribution. Meanwhile, we randomly generate an integer, indicating the object label, within the range of the total category number of OmniObject3D. The integer is then transformed into a label embedding by an embedding layer. This enables the model to be trained on different taxonomies simultaneously. Next, we forward the label embedding and latent code $\mathbf{z}$ into a mapping network and generate an intermediate latent code $\mathbf{w}$ $(1 \times 15 \times 512)$ with promoted dimensions. Here, the first 14 dimensions of $\mathbf{w}$ ($\mathbf{w}_1$ in the figure) are used as the input to the object generator ($\mathbf{G}_1$ in the figure) while the last dimension of $\mathbf{w}$ ($\mathbf{w}_2$ in the figure) are employed as the input to the color generator ($\mathbf{G}_2$ in the figure). $\mathbf{G}_1$ and $\mathbf{G}_2$ generate the tri-plane representation $(3 \times 32 \times 64 \times 64)$ and the color tensor $(3 \times 10)$, respectively. Following [18], we use the color generator which allows convenient texture conversion. After that, we determine queries, $\mathbf{Q}$, based on the randomly generated camera pose. The tri-plane representation is then taken as keys, $\mathbf{K}$, and transmitted into the tri-plane decoder together with $\mathbf{Q}$. Then, we use a decoder to generate the SDF of the object, which can be sampled to create a 3D mesh. This decoder also produces a feature as the input to a view direction mapper, from which the output is multiplied by the color tensor, $\mathbf{V}$ (the values), in the

Figure 2. Visualization results of the synthesis images rendered from the NeRFs generated by our solution.



Figure 3. Visualization results of the synthesis images rendered from the NeRFs generated by our solution without additional label embedding.

| Method | FID |
|---|---|
| Δ= Ours w/ label embedding - Ours w/o label embedding | -4.8 |

Table 1. Improvement in terms of FID of our solution after using the additional label embedding.

attention module. Finally, the attention module produces the 2D synthesis image and a discriminator [2] which only functions during training is used for adversarial training.

**Model Training.** Since we are simultaneous optimizing NeRFs and SDFs representations, we ensure the generator can initially generate a unit sphere by initializing our model with SDF-pretraining as [18]. Then, the generator is optimized using the Adam optimizer with a learning rate of 0.0025, while the discriminator [2] uses a learning rate of 0.002. The batch size is set to 32 and the model is trained for 300,000 iterations. Additionally, we also use adaptive discriminator augmentation [8] to improve model generalization. Overall, the model training objective is formulated as:

$$\alpha_0 \mathcal{L}_{path} + \alpha_1 \mathcal{L}_e + \alpha_2 \mathcal{L}_v + \alpha_3 \mathcal{L}_{sdf}, \quad (1)$$

where $\mathcal{L}_{path}$ denotes the Path length regulation loss [10], $\mathcal{L}_e$ and $\mathcal{L}_v$ denote the entropy loss and total variation loss [10], respectively, $\mathcal{L}_{sdf}$ represents the SDF eikonal loss [18], and $\alpha_i$ are weighting parameters. We set $\alpha_0$=2, $\alpha_1$=0.05, $\alpha_2$=0.5 and $\alpha_3$=0.1.

## 3. Experiments

**Dataset.** We train our model from scratch for 3D object generation on the training dataset provided by the ICCV 2023 OmniObject3D Challenge. The complete OmniObject3D dataset [23] contains 216 classes from a wide range of daily categories and approximately 6k objects in total where each object has 100 images in different views.

**Efficient Model Training.** Instead of training one model for each class, we train a single model with various types of objects. The model takes each image along with its camera pose, focal length and object type. Besides, we empirically find that we can effectively train this model with only a few images per object from various classes, so we select the first 8 images of each object for model training. This significantly improves the model training efficiency.

**Results.** Fig. 2 shows some visualization results of synthesis images generated by our solution on OmniObject3D. It can be seen that our approach is capable of generating high-fidelity images from different views for 3D object generation. Besides, Fig. 3 shows results of our solution without the additional label embedding. Comparing Fig. 2 with Fig. 3, we can see that without the additional label embedding, the quality of the synthesis images is significantly worse than the original solution. Furthermore, Table. 1 shows that using the additional label embedding can improve the solution by 4.86 in terms of FID. This verifies that the additional label embedding is beneficial for improving model generalization.

## 4. Conclusion

In this technical report, we present a solution for 3D object detection of ICCV 2023 OmniObject3D Challenge. The key insight of our solution is to learn effective NeRFs and SDFs Representations with 3D generative adversarial networks for 3D object generation. Overall, this solution is one of the final top-3-place solutions in the ICCV 2023 OmniObject3D Challenge.

## References

[1] Heli Ben-Hamu, Haggai Maron, Itay Kezurer, Gal Avineri, and Yaron Lipman. Multi-chart generative surface modeling.

*ACM Transactions on Graphics*, 37(6):1–15, 2018. 1

[2] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 1, 2, 3

[3] Jaehoon Choi, Dongki Jung, Taejae Lee, Sangwook Kim, Youngdong Jung, Dinesh Manocha, and Donghwan Lee. Tmo: Textured mesh acquisition of objects with a mobile device by using differentiable rendering. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16674–16684, 2023. 1

[4] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021. 1

[5] Matheus Gadelha, Subhransu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. In *2017 International Conference on 3D Vision (3DV)*, pages 402–411. IEEE, 2017. 1

[6] Xintong Han, Xiaojun Hu, Weilin Huang, and Matthew R Scott. Clothflow: A flow-based model for clothed person generation. In *Proc. of the IEEE/CVF international conference on computer vision*, pages 10471–10480, 2019. 1

[7] Lutao Jiang, Ruyi Ji, and Libo Zhang. Sdf-3dgan: A 3d object generative method based on implicit signed distance function. *arXiv preprint arXiv:2303.06821*, 2023. 1

[8] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *Advances in neural information processing systems*, 33:12104–12114, 2020. 3

[9] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 1

[10] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proc. of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 1, 2, 3

[11] Yibo Liu, Kelly Zhu, Guile Wu, Yuan Ren, Bingbing Liu, Yang Liu, and Shan Jinjun. Mv-deepsdf: Implicit modeling with multi-sweep point clouds for 3d vehicle reconstruction in autonomous driving. In *Proc. of the IEEE/CVF International Conference on Computer Vision*, 2023. 1

[12] Sebastian Lunz, Yingzhen Li, Andrew Fitzgibbon, and Nate Kushman. Inverse graphics gan: Learning to generate 3d shapes from unstructured 2d data. *arXiv preprint arXiv:2002.12674*, 2020. 1

[13] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proc. of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019. 1

[14] 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. *Communications of the ACM*, 65(1):99–106, 2021. 1

[15] Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas J Guibas. Structurenet: Hierarchical graph networks for 3d shape generation. *arXiv preprint arXiv:1908.00575*, 2019. 1

[16] Norman Müller, Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Bulo, Peter Kontschieder, and Matthias Nießner. Diffrf: Rendering-guided 3d radiance field diffusion. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4328–4338, 2023. 1, 2

[17] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proc. of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019. 1

[18] Dario Pavllo, David Joseph Tan, Marie-Julie Rakotosaona, and Federico Tombari. Shape, pose, and appearance from a single image via bootstrapped radiance field inversion. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4391–4401, 2023. 1, 2, 3

[19] Mathis Petrovich, Michael J Black, and Gül Varol. Action-conditioned 3d human motion synthesis with transformer vae. In *Proc. of the IEEE/CVF International Conference on Computer Vision*, pages 10985–10995, 2021. 1

[20] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proc. of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 1

[21] Tengfei Wang, Bo Zhang, Ting Zhang, Shuyang Gu, Jianmin Bao, Tadas Baltrusaitis, Jingjing Shen, Dong Chen, Fang Wen, Qifeng Chen, et al. Rodin: A generative model for sculpting 3d digital avatars using diffusion. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4563–4573, 2023. 1, 2

[22] Xiaogang Wang, Marcelo H Ang, and Gim Hee Lee. Voxel-based network for shape completion by leveraging edge generation. In *Proc. of the IEEE/CVF international conference on computer vision*, pages 13189–13198, 2021. 1

[23] Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Jiawei Ren, Liang Pan, Wayne Wu, Lei Yang, Jiaqi Wang, Chen Qian, et al. Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 803–814, 2023. 1, 2, 3