
VPP: Efficient Conditional 3D Generation via Voxel-Point Progressive Representation

Supplementary Material

Anonymous Author(s)

Affiliation

Address

email

A Implementation Details

A.1 Experimental Details

Training Details We use ShapeNetCore from ShapeNet [1] as the training dataset. ShapeNet is a clean set of 3D CAD object models with rich annotations, including 51K unique 3D models from 55 common object categories. For the acquisition of multi-modal data, we follow ReCon [8] for multi-view rendering and utilize BLIP [4] based on the rendered images to obtain textual data. In Table 1, we show the training hyperparameters and model architecture information of each part of our VPP.

Config	3D VQGAN	Voxel Generator	Grid Smoother	Point Upsampler
Training Parameters				
Optimizer	Adam	AdamW	AdamW	AdamW
Learning rate	1e-4	1e-3	1e-3	1e-3
Weight decay	1e-4	5e-2	5e-2	5e-2
Training epochs	300	100	100	300
Learning rate scheduler	cosine	cosine	cosine	cosine
Batch size	128	128	128	128
Drop path rate	-	0.1	0.1	0.1
Input point size	8192	8192	8192	1024
Model Architecture				
Backbone	CNN	Transformer	Transformer	Transformer
Layers	6	12	4	6
Hidden size	384	384	64	384
Heads	-	6	4	6
Voxel resolution	24	24	24	24
Point patch size	-	-	-	32
GPU device	NVIDIA A100	NVIDIA A100	NVIDIA A100	NVIDIA A100

Table 1: Training recipes for 3D VQGAN, Voxel Generator, Grid Smoother and Point Upsampler.

Downstream Tasks Details Following Point-E [6], we use pre-trained PointNet++ as a classifier in all ACC, FID, and IS evaluations to extract the features and calculate the accuracy of generated point clouds. In point cloud generation and editing, we employ 8 or 4 steps for a parallel generation. The generation task utilizes initial voxel codebooks composed entirely of [MASK] tokens. In editing, we

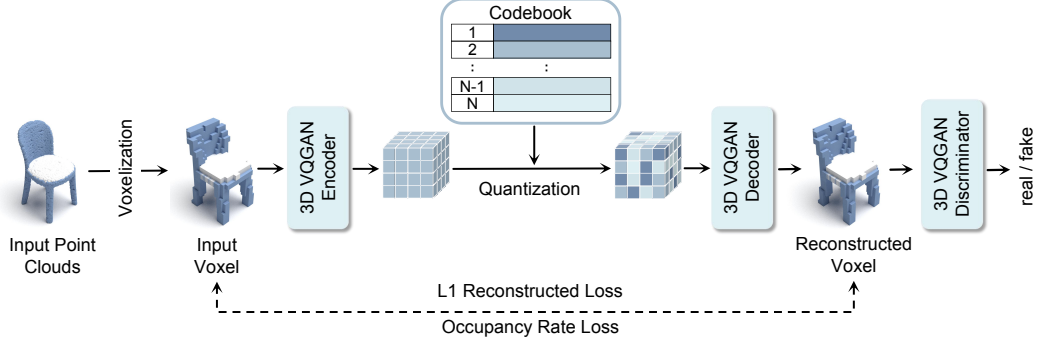


Figure 1: The training overview of 3D VQGAN. We introduce the occupancy rate loss to have a better reconstruction of the voxel.

13 extract VQGAN features from the original point cloud to initialize the voxel codebooks. As for the
 14 transfer classification task on ScanObjectNN [10] and ModelNet40 [14], we fully follow the previous
 15 work [2, 7] configuration and trained 300 epochs with the AdamW optimizer, and used the voting
 16 strategy in testing.

17 A.2 Implementation Details of 3D VQGAN

18 We show the detailed training overview of 3D VQGAN in Figure 1. Following the traditional training
 19 of VQGAN [3], we use L1 loss to supervise the reconstruction of voxel, and feed the reconstructed
 20 voxel into the discriminator to ensure the generated authenticity by GAN loss. Besides the L1 loss and
 21 GAN loss, we also introduce the occupancy rate loss to make the occupancy rate of the reconstructed
 22 voxel grid similar to the ground truth, so as to obtain a better reconstruction of the voxel.

23 B Additional Experiments

24 We conduct more experiments to further demonstrate the generation quality and universality of VPP.
 25 Including diversity & specific text-conditional generation, point-based mesh reconstruction, and
 26 few-shot transfer classification.

27 B.1 Diversity & Specific Text-Conditional Generation

28 We show the diversity qualitative results of VPP on text-conditional 3D generation in Figure 2 (a).
 29 It can observe that VPP can generate a broad category of shapes with rich diversity while remain
 30 faithful to the provided text descriptions. Meanwhile, we present the qualitative results of VPP on
 31 more specific text-conditional 3D generation in Figure 2 (b). Notably, VPP can generate high-fidelity
 32 shapes that well react to very specific text descriptions, like 'a round chair with arms', 'a
 33 round table with four legs' and 'an old-style propeller aircraft', etc.

34 B.2 Point-Based Mesh Reconstruction

35 Besides the representation of point clouds, we also show the common mesh representation of
 36 generated shapes to demonstrate the high generation quality of VPP. We use DMTet [9] to reconstruct
 37 the mesh from the generated point clouds, where the reconstructed examples are presented in Figure 3.
 38 It is observed that the reconstructed mesh is still high-quality with fine geometric details, which
 39 further proves that VPP not only has high fidelity generation but also supports the output of mesh
 40 representation.

41 B.3 Few-Shot Transfer Classification

42 We conduct few-shot experiments on the ModelNet40 [14] dataset, and the results are shown in
 43 Table 2. In the self-supervised benchmark without the use of additional modality data, VPP achieves
 44 excellent performance compared to previous works.

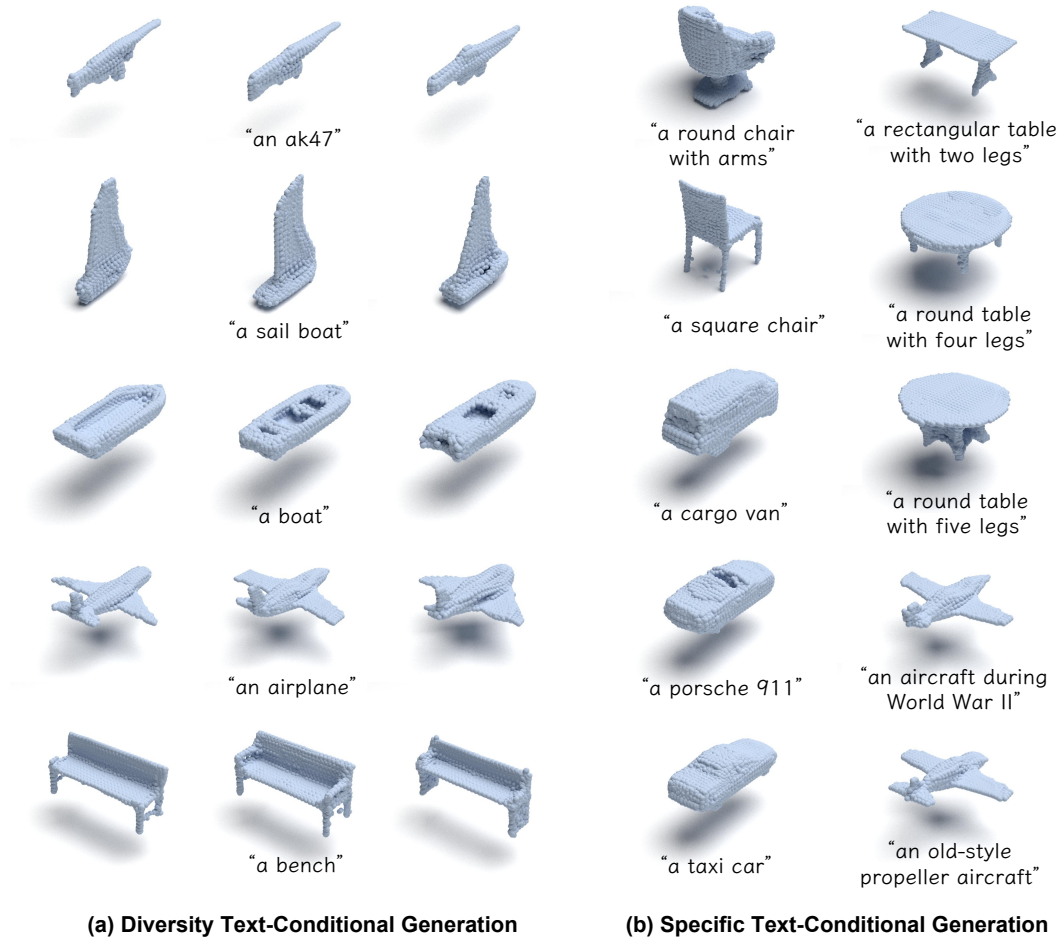


Figure 2: (a) **Diversity** qualitative results of VPP on text-conditional 3D generation. (b) Qualitative results of VPP on more **specific** text-conditional 3D generation. VPP can generate a broad category of shapes with rich diversity and high fidelity, while remain faithful to the provided text descriptions. Besides, VPP can react to very specific text descriptions, like "a round chair with arms".

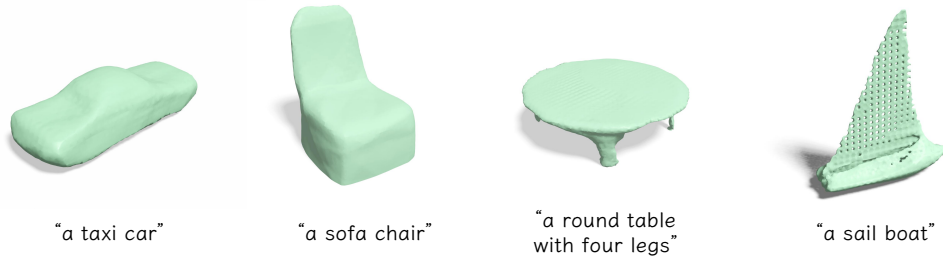


Figure 3: Examples of reconstructed mesh from the generated point clouds. VPP not only has high-fidelity generation but also supports the output of mesh representation.

Method	5-way		10-way	
	10-shot	20-shot	10-shot	20-shot
<i>Supervised Learning Only</i>				
DGCNN [13]	31.6 \pm 2.8	40.8 \pm 4.6	19.9 \pm 2.1	16.9 \pm 1.5
OcCo [12]	90.6 \pm 2.8	92.5 \pm 1.9	82.9 \pm 1.3	86.5 \pm 2.2
<i>with Self-Supervised Representation Learning</i>				
Transformer [11]	87.8 \pm 5.2	93.3 \pm 4.3	84.6 \pm 5.5	89.4 \pm 6.3
OcCo [12]	94.0 \pm 3.6	95.9 \pm 2.3	89.4 \pm 5.1	92.4 \pm 4.6
Point-BERT [15]	94.6 \pm 3.1	96.3 \pm 2.7	91.0 \pm 5.4	92.7 \pm 5.1
MaskPoint [5]	95.0 \pm 3.7	97.2 \pm 1.7	91.4 \pm 4.0	93.4 \pm 3.5
Point-MAE [7]	96.3 \pm 2.5	97.8 \pm 1.8	92.6 \pm 4.1	95.0 \pm 3.0
Point-M2AE [16]	96.8 \pm 1.8	98.3 \pm 1.4	92.3 \pm 4.5	95.0 \pm 3.0
VPP (ours)	96.9 \pm 1.9	98.3 \pm 1.5	93.0 \pm 4.0	95.4 \pm 3.1

Table 2: Few-shot transfer classification results on ModelNet40. Overall accuracy (%) without voting is reported.

46 C Limitation and Future Work

47 Although, as a generative model, the proposed VPP has achieved the unified model on multiple
48 downstream 3D tasks, the current work focuses on some relatively simple tasks. For instance, we
49 perform 3D completion task on point cloud upsampling, and the pre-training transfer learning on
50 3D classification. It would be more valuable to extend it to more complex 3D applications, e.g.,
51 shape completion and 3D pre-training on detection & segmentation. Furthermore, the current model
52 is trained with ShapeNet [1], which is still a limited-scale dataset across different categories. It
53 would be promising to explore large-scale 3D dataset for training the model to further improve the
54 generation quality and capability.

55 D Broader Impact

56 Conditional 3D generation has the potential to improve the 3D design and assist practitioners
57 efficiently to work on content creation. Similarly, we propose VPP with the hope of enhancing better
58 creativity and work as an alternative tool for artists to more efficiently design. Besides, conditional
59 3D generation methods allow the public to have more convenient access to 3D craft, which can bring
60 many benefits to human-AI collaborative application.

61 References

- 62 [1] Angel X. Chang, Thomas A. Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qi-Xing Huang, Zimo Li,
63 Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet:
64 An information-rich 3d model repository. *CoRR*, abs/1512.03012, 2015. 1, 4
- 65 [2] Runpei Dong, Zekun Qi, Linfeng Zhang, Junbo Zhang, Jianjian Sun, Zheng Ge, Li Yi, and Kaisheng
66 Ma. Autoencoders as cross-modal teachers: Can pretrained 2d image transformers help 3d representation
67 learning? In *Int. Conf. Learn. Represent. (ICLR)*, 2023. 2
- 68 [3] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image
69 synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages
70 12873–12883, 2021. 2
- 71 [4] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. BLIP: bootstrapping language-image pre-
72 training for unified vision-language understanding and generation. In *ICML*, volume 162 of *Proceedings*
73 *of Machine Learning Research*, pages 12888–12900. PMLR, 2022. 1
- 74 [5] Haotian Liu, Mu Cai, and Yong Jae Lee. Masked discrimination for self-supervised learning on point
75 clouds. In *Eur. Conf. Comput. Vis. (ECCV)*, 2022. 4
- 76 [6] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for
77 generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022. 1

- 78 [7] Yatian Pang, Wenxiao Wang, Francis E. H. Tay, Wei Liu, Yonghong Tian, and Li Yuan. Masked autoen-
79 coders for point cloud self-supervised learning. In *Eur. Conf. Comput. Vis. (ECCV)*, 2022. 2, 4
- 80 [8] Zekun Qi, Runpei Dong, Guofan Fan, Zheng Ge, Xiangyu Zhang, Kaisheng Ma, and Li Yi. Contrast with
81 reconstruct: Contrastive 3d representation learning guided by generative pretraining. In *Proc. Int. Conf.*
82 *Mach. Learn. (ICML)*, 2023. 1
- 83 [9] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: a
84 hybrid representation for high-resolution 3d shape synthesis. In *NeurIPS*, pages 6087–6101, 2021. 2
- 85 [10] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Thanh Nguyen, and Sai-Kit Yeung. Revisiting
86 point cloud classification: A new benchmark dataset and classification model on real-world data. In
87 *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 1588–1597, 2019. 2
- 88 [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
89 Kaiser, and Illia Polosukhin. Attention is all you need. In *Adv. Neural Inform. Process. Syst. (NIPS)*, pages
90 5998–6008, 2017. 4
- 91 [12] Hanchen Wang, Qi Liu, Xiangyu Yue, Joan Lasenby, and Matt J Kusner. Unsupervised point cloud
92 pre-training via occlusion completion. In *Proceedings of the IEEE/CVF international conference on*
93 *computer vision*, pages 9782–9792, 2021. 4
- 94 [13] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon.
95 Dynamic graph CNN for learning on point clouds. *ACM Trans. Graph.*, 38(5):146:1–146:12, 2019. 4
- 96 [14] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao.
97 3d shapenets: A deep representation for volumetric shapes. In *IEEE/CVF Conf. Comput. Vis. Pattern*
98 *Recog. (CVPR)*, pages 1912–1920, 2015. 2
- 99 [15] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert: Pre-training 3d
100 point cloud transformers with masked point modeling. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog.*
101 *(CVPR)*, 2022. 4
- 102 [16] Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hongsheng Li.
103 Point-m2AE: Multi-scale masked autoencoders for hierarchical point cloud pre-training. In *Adv. Neural*
104 *Inform. Process. Syst. (NeurIPS)*, 2022. 4