

3 2 More Examples of Shape-Conditioned Multimodal Generation

4 In Figure 3 and Figure 4, we showcase more examples of point cloud captioning and point cloud-
5 conditioned image generation.



Figure 3: **Point cloud captioning.** In each row, we show the input point clouds on the left and the generated captions on the right.

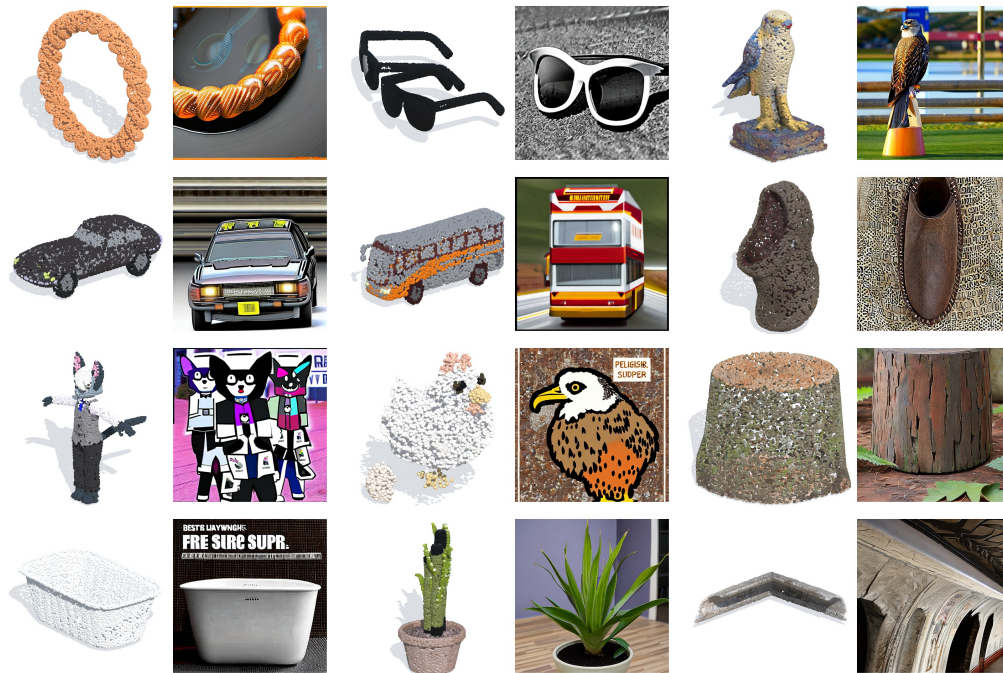


Figure 4: **Point cloud-conditioned image generation.** Each row shows three examples (input point clouds and generated images).

6 3 Details on Raw Text Generation and Filtering

7 3.1 Raw Text Generation

8 We leverage the metadata from the four datasets to generate the raw texts. Although the original
9 datasets may contain numerous attributes for each shape, we carefully choose the most informative
10 ones to compose the text, ensuring its quality and relevance.

11 **Objaverse:** We utilize the name associated with each shape to serve as the text.

12 **ShapeNetCore:** For each shape, we generate three types of texts: (a) the name, (b) the category
13 name (with a total of 55 categories), and (c) the concatenation of the sub-category names (with a
14 total of 336 sub-categories), separated by commas.

15 **3DFuture:** For each shape, we generate two types of texts: (a) the category, and (b) the concatena-
16 tion of category, style, theme, and material, separated by commas.

17 **ABO:** For each shape, we generate two types of texts: (a) the `item_name`, and (b) the `product_type`.

18 In this way, we generate one or more raw texts for each shape.

19 3.2 Raw Text Filtering

20 We employ GPT-4 [4] to filter out uninformative raw texts. To accomplish this, we divide all the raw
21 texts into batches, each containing 256 entries, and process each batch independently using GPT-4.
22 Here is an example illustrating the prompt we used and the corresponding response generated by
23 GPT-4.

```
I am analyzing a 3D dataset with various text descriptions for the 3D models.
However, many of these texts are inaccurate or uninformative, and therefore,
not suitable as descriptions for 3D models. I need your help to identify such
incorrect texts. Specifically, if a text primarily consists of irrelevant or unin-
formative content, such as timestamps, model numbers, incomprehensible
descriptions, random filenames (e.g., "my project"), random characters, etc.,
please respond with "N". If a text contains a clear noun (or noun phrase)
that could potentially describe a 3D object, please respond with "Y". You
will find a list of texts below, and each line contains a three-digit ID and
associated text. For each text, please respond with "Y" or "N", following the
ID number (e.g., "001 Y" or "002 N"). Please evaluate all 256 texts.
```

```
000 New project ( 19 )
001 3December - Chemistry
002 Fake Brand Soda Can
003 Spartan Shild
004 Apple3d
005 Landmine
006 FaunveinB-S
007 FIGURA 5
008 Sphero Blue
009 Sofa
010 Maddox
011 A3 Complete
012 Suspension Bridge
013 Maung
014 Captain-americas-shield
015 sphorb4
.....
```

```
000 N
001 Y
002 Y
003 Y
004 Y
005 Y
006 N
007 N
008 Y
009 Y
010 N
011 N
012 Y
013 N
014 Y
015 N
.....
```

24
25 Afterward, we combine all the responses to create the final filtering results, effectively removing
26 approximately 30% of the raw texts.

27 4 Details on the Backbone Scaling Experiment

28 In Figure 4 of the main paper, we investigate the performance and scalability of various backbones
29 when scaling up their model sizes. For this experiment, we employ a default resolution of 10,000
30 points for input point clouds, a batch size of 200, and conduct the experiment on a single A100 GPU.
31 In general, if instructions are given in the original paper of a backbone, we scale up the model as
32 instructed. Otherwise, we scale up the model by expanding width or depth (i.e., stacking blocks or
33 layers). Specifically, we scale up each backbone as follow:

34 **PointBERT [13]** The scaling parameters are shown in Table 1. We scaled PointBERT to 72.1M
35 parameters beyond the 32.3M version reported in Figure 4 of the main paper. However, at this scale,
36 the model dramatically overfits on the training data and performs worse on all benchmarks than the
37 32.3M version.

Table 1: Hyperparameters for scaling up PointBERT [13].

# Parameters	# Layers	Width	# Heads	MLP Dim	# Patches	Patch Embed Dim
5.1M	6	256	4	1024	64	96
13.3M	6	512	8	1024	64	128
32.3M	12	512	8	1536	384	256
72.1M	12	768	12	2304	512	256

38 **SparseConv [1]** The smallest version (5.3M parameters) of the model is adapted from the
39 MinkowskiFCNN model by adjusting the width of the final convolution and linear layers. The
40 remaining three models are adaptations of MinkowskiResNet, each varying in the number of basic
41 ResNet blocks used. See Table 2 for the specific scaling parameters.

Table 2: Hyperparameters for scaling up SparseConv [1].

# Parameters	# Convolution Layers	# Linear Layers
5.3M	7	4
29.0M	18	3
33.7M	26	3
41.3M	42	3

42 **PointNeXt [7]** PointNeXt is proposed as a scalable version of PointNet++ [6], and includes
43 S/B/L/XL variants in the original paper. We simply adopt these official configurations.

44 **DGCNN [10] and PointNet [5]** For these two backbones without a hierarchical structure, we
45 increase the width of each layer proportionally to scale up to 4xPointNet and 2xDGCNN before we
46 hit the GPU memory limit. As the models operate completely on dense points, it is impractical to use
47 the default 10k-point resolution. We thus reduce the input resolution for the two backbones, resulting
48 in 1k points for DGCNN and 4k points for PointNet.

49 5 Details on Training and Evaluation

50 **Training Details** We freeze the CLIP text and image encoders and train the 3D encoder and two
51 projection heads on our ensembled dataset using the cross-modal contrastive loss. We train the
52 model on a single A100 GPU with a batch size of 200. Since we precache the text and image CLIP
53 embeddings of all shapes, the training is greatly accelerated and takes about 300 A100 hours for
54 convergence. We utilize an exponential learning rate schedule, and employ a range test to find the
55 initial learning rate. For 32.3M version of PointBERT, we utilize a learning rate of $5e - 4$; for 72.1M
56 version of PointBERT, we utilize a learning rate of $4e - 4$; and for other models, we utilize a learning
57 rate of $1e - 3$. For hard-negative mining, the number of seed shapes s is set to 40, and the number of
58 neighbors m is set to 5 per shape, and the threshold δ is set to 0.1.

59 **Fine-tuning CLIP Text and Image Encoders?** After training OpenShape-PointBERT, we con-
60 ducted experiments to unfreeze and finetune the CLIP text encoder for a single epoch. However, the
61 results obtained did not demonstrate any noticeable improvement on the benchmarks. Moreover,
62 we observed that finetuning the CLIP text encoder could potentially undermine the generalization
63 capabilities of CLIP and hinder the integration of OpenShape embeddings into existing CLIP-based
64 models. As a result, we choose to freeze the CLIP encoders throughout the entire training process.

65 **Evaluation Details** We evaluated all baselines using their publicly released pretrained checkpoints.
66 Additionally, we retrained ULIP [12] on our ensembled training shapes using their official code base
67 and backbone networks. Note that the retrained ULIP model utilized the original raw texts from the
68 four datasets during training (prompt engineering is also applied), rather than our filtered and enriched
69 texts. For ModelNet40 [11], the evaluation is conducted on the test split with 2,468 shapes. Regarding
70 ScanObjectNN [9], we follow ULIP [12] to evaluate on the OBJ_ONLY version, which contains
71 581 test shapes. For Objaverse-LVIS [2], the input is 10,000 sampled points with point colors. For
72 ModelNet40 [11], the input is 10,000 sampled points without color. For ScanObjectNN [9], we utilize
73 the official 2,048 points without color as input. All methods use the same input during evaluation.
74 The forward inference time on an A100 GPU for a 10,000-point point cloud is approximately 0.9ms
75 for OpenShape-SparseConv and 3.8ms for OpenShape-PointBERT.

76 6 Details on Shape-Conditioned Multimodal Generation

77 **Point Cloud Captioning** CLIPCap [3] utilizes a 10-token prefix generated from CLIP image
78 embeddings to enable GPT-2 for captioning. In order to align with the off-the-shelf CLIPCap model,
79 we trained a variant of OpenShape-PointBERT that employs CLIP ViT-B/32 embeddings instead
80 of OpenCLIP ViT-G/14 used in other experiments. Consequently, we directly input the point cloud
81 encoding, *without normalization*, into CLIPCap for captioning.

82 **Point Cloud Conditioned Image Generation** We take the Stable Diffusion v2.1 unCLIP model [8]
83 for image generation and replace the CLIP image condition encoder with our OpenShape encoder to
84 perform image generation conditioned on point clouds (and optionally text prompts). The unCLIP
85 model takes CLIP ViT-L/14 embeddings without normalization as input. To match the embedding
86 space, we trained a variant of OpenShape-PointBERT with CLIP ViT-L/14 embeddings. Additionally,
87 we noticed a significant mismatching of scales (L_2 -norm of embedding vectors) between ViT-L/14
88 image embeddings and OpenShape embeddings. To mitigate this issue, we perform a re-normalization
89 on OpenShape embeddings to a L_2 -norm of $\frac{1}{2}\sqrt{768}$, which is our observed mean L_2 -norm of ViT-
90 L/14 image embeddings. We use 50 diffusion steps. The guidance scale can be tuned freely.

91 References

- 92 [1] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets:
93 Minkowski convolutional neural networks. In *Proceedings of the IEEE Conference on Computer
94 Vision and Pattern Recognition*, pages 3075–3084, 2019.
- 95 [2] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt,
96 Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe
97 of annotated 3d objects. *arXiv preprint arXiv:2212.08051*, 2022.
- 98 [3] Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning.
99 *arXiv preprint arXiv:2111.09734*, 2021.
- 100 [4] OpenAI. Gpt-4 technical report, 2023.
- 101 [5] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point
102 sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer
103 vision and pattern recognition*, pages 652–660, 2017.
- 104 [6] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical
105 feature learning on point sets in a metric space. *Advances in neural information processing
106 systems*, 30, 2017.
- 107 [7] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Hammoud, Mohamed Elhoseiny,
108 and Bernard Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling
109 strategies. *arXiv:2206.04670*, 2022.

- 110 [8] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
111 text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 112 [9] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Thanh Nguyen, and Sai-Kit Yeung.
113 Revisiting point cloud classification: A new benchmark dataset and classification model on
114 real-world data. In *Proceedings of the IEEE/CVF international conference on computer vision*,
115 pages 1588–1597, 2019.
- 116 [10] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M
117 Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics*
118 (*tog*), 38(5):1–12, 2019.
- 119 [11] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and
120 Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of*
121 *the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015.
- 122 [12] Le Xue, Mingfei Gao, Chen Xing, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu,
123 Juan Carlos Niebles, and Silvio Savarese. Ulip: Learning unified representation of language,
124 image and point cloud for 3d understanding. *arXiv preprint arXiv:2212.05171*, 2022.
- 125 [13] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert:
126 Pre-training 3d point cloud transformers with masked point modeling. In *Proceedings of the*
127 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.