

MASH: Masked Anchored SpHerical Distances for 3D Shape Representation and Generation - Supplemental Material

CHANGHAO LI, University of Science and Technology of China, China

YU XIN, University of Science and Technology of China, China

XIAOWEI ZHOU, State Key Laboratory of CAD & CG, Zhejiang University, China

ARIEL SHAMIR, Reichman University, Israel

HAO ZHANG, Simon Fraser University, Canada

LIGANG LIU, University of Science and Technology of China, China

RUIZHEN HU, Shenzhen University, China

CCS Concepts: • Computing methodologies → Shape modeling; Neural networks.

ACM Reference Format:

Changhao Li, Yu Xin, Xiaowei Zhou, Ariel Shamir, Hao Zhang, Ligang Liu, and Ruizhen Hu. 2025. MASH: Masked Anchored SpHerical Distances for 3D Shape Representation and Generation - Supplemental Material . *ACM Trans. Graph.* 1, 1 (April 2025), 8 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 SHAPE APPROXIMATION

Optimization details. To get more uniform sampling on the surfaces determined by the anchors, we set the number of uniformly pre-sampled directions $N_{\text{dir}} = 1000$, and then keep 80% of sampled points of all anchors by the farthest point sampling algorithm. Besides, in order to calculate the boundary-continuous loss term, we sample another $N_{\text{mask}} = 90$ points on the vision mask boundary of each anchor. When optimizing MASH parameters, we use AdamW [Loshchilov 2017] as our optimizer with an initial learning rate of 2e-3, and we use ReduceLROnPlateau as our scheduler with a final learning rate of 1e-3, a factor of 0.8, and patience of 2.

Hyperparameter choices. The key hyperparameters in our MASH representation are anchor numbers M , mask degrees K , and spherical harmonic degrees L . When L is 0 or 1, the surface patch of each anchor approximates a plane or a near-spherical surface, and when L exceeds 2, the excessive number of SH basis functions leads to increased computational costs. Therefore, to balance accuracy and computational efficiency, we choose L to be 2. After extensive experiments, we find that regarding the object complexity of the

Authors' addresses: Changhao Li, University of Science and Technology of China, Hefei, China, lch0510@mail.ustc.edu.cn; Yu Xin, University of Science and Technology of China, Hefei, China, xy0731@mail.ustc.edu.cn; Xiaowei Zhou, State Key Laboratory of CAD & CG, Zhejiang University, Hangzhou, China, xwzhou@zju.edu.cn; Ariel Shamir, Reichman University, Herzliya, Israel, arik@runi.ac.il; Hao Zhang, Simon Fraser University, Vancouver, Canada, haoz@cs.sfu.ca; Ligang Liu, University of Science and Technology of China, Hefei, China, lglu@ustc.edu.cn; Ruizhen Hu, Shenzhen University, Shenzhen, China, ruižhen.hu@gmail.com.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Association for Computing Machinery.

0730-0301/2025/4-ART \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Table 1. Quantitative comparison of shape approximation results with different mask and SH degrees. Here we set the number of MASH anchors to be 400 for fair comparison.

	K=4	K=3	K=2	K=1	K=0
Chamfer ↓	5.434	5.450	5.698	6.392	11.281
F-Score ↑	0.997	0.997	0.996	0.993	0.872
	L=4	L=3	L=2	L=1	L=0
Chamfer ↓	5.323	5.411	5.450	6.828	12.593
F-Score ↑	0.998	0.997	0.997	0.989	0.814

Table 2. Quantitative comparison of shape approximation results with different numbers of MASH anchors, where $K = 3$ and $L = 2$.

	M=400	M=200	M=100	M=50	M=20	M=10
Chamfer ↓	5.450	6.013	7.174	11.267	16.536	18.990
F-Score ↑	0.997	0.994	0.988	0.924	0.851	0.812
Time (s)	39.027	28.560	23.979	20.817	18.083	14.352

Table 3. Quantitative comparison of results with and without the inverse transform during the MASH optimization process.

	Chamfer ↓	F-Score ↑
with Inversion	5.450	0.997
w/o Inversion	5.834	0.982

ShapeNet-V2 dataset, setting the mask degree $K = 3$ and the anchor number $M = 400$ can accurately represent most of the objects in this dataset while maintaining high computational efficiency. Table 1 shows how the approximation error changes with different K and L when $M = 400$. Table 2 shows how the approximation error changes with different M with $K = 3$ and $L = 2$.

Importance of Inverse Transformation. We also conduct an experiment to show the importance of the inverse transformation. If inverse transformation is not used, the representation ability of each anchor will decrease and the planar area of the given shape will become one of the most difficult parts to be fit with spherical harmonics. The results in Table 3 confirmed the effectiveness of applying the inverse transformation.

Table 4. Comparison of computation time and memory usage when different numbers of points Q are sampled on the given shape, under the default setting of $M=400$, $K=3$, and $L=2$.

$ Q $	PS	L_f	L_c	L_ω	TC	Memory	Chamfer
1K	2.806	3.372	3.391	3.416	13.236	1690	19.901
2K	4.307	4.421	4.386	4.129	17.628	1691	13.684
5K	4.336	4.878	4.797	4.651	18.852	1703	11.257
10K	5.692	7.136	7.210	7.205	27.104	1719	6.833
20K	5.727	7.902	7.897	7.884	29.371	1745	5.263
40K	7.675	10.366	10.263	9.901	38.762	1839	4.971

Timing and memory usage. Table 4 shows the detailed computation time and memory usage during the MASH optimization process when different numbers of points are sampled from the given shape, averaging over 100 shapes sampled from the ShapeNet-V2 dataset. The units for time are the second, and the units for memory usage are MB. Specifically, Point Sampling (PS) represents the total time taken to sample points from MASH, L_f , L_c , and L_ω correspond to the total time taken for the corresponding loss calculations, and Time Consumption (TC) is the total duration of the optimization algorithm.

Note that when fitting point clouds with different sizes, we need to correspondingly adjust the number of pre-sampled directions N_{dir} to prevent overfitting. Therefore, we set $N_{\text{dir}} = 10 \cdot |Q|/M$. As $|Q|$ increases, anchors can more accurately expand and cover along the object's surface, thereby gradually improving the fitting efficiency.

Scene MASH optimization. Other than 3D shapes, we also make an additional attempt to fit the point cloud of the indoor scenes with MASH by increasing the anchor number M to 800. Some example results are shown in Figure 1, which demonstrates the potential of MASH in the representation and processing of large-scale data.

2 SURFACE RECONSTRUCTION

Reconstruction details. Most of the technical details about the surface reconstruction have been provided in the main paper. Here we show how the normal n_{scr} of the sampled points is computed:

$$n_x = \frac{\partial y}{\partial \phi} \frac{\partial z}{\partial \theta} - \frac{\partial z}{\partial \phi} \frac{\partial y}{\partial \theta} \quad (1)$$

$$n_y = \frac{\partial z}{\partial \phi} \frac{\partial x}{\partial \theta} - \frac{\partial x}{\partial \phi} \frac{\partial z}{\partial \theta} \quad (2)$$

$$n_z = \frac{\partial x}{\partial \phi} \frac{\partial y}{\partial \theta} - \frac{\partial y}{\partial \phi} \frac{\partial x}{\partial \theta} \quad (3)$$

where $\{\theta, \phi\}$ are the spherical coordinates and $\{x, y, z\}$ is the position of that sampled point.

Evaluation metrics. We compare our reconstruction results with ground truth (GT) meshes using common evaluation metrics, including *Hausdorff distance* (D_H), *L1 Chamfer Distance* (CD), *F-Score* and *mesh cosine similarity* (S_{cos}) commonly used in previous works [Erler et al. 2020; Fan et al. 2017; Huttenlocher et al. 1993; Lin et al. 2022; Park et al. 2019], which are defined as:

$$D_H(P, Q) = \max(\max_{p \in P} d(p, Q), \max_{q \in Q} d(q, P)) \quad (4)$$

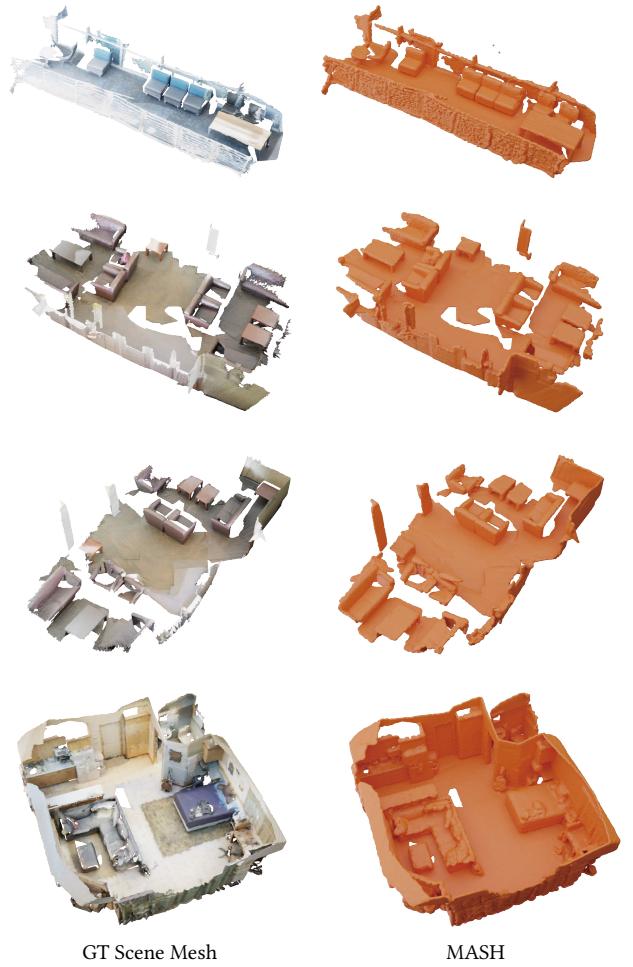


Fig. 1. Qualitative results on fitting and reconstruction for large-scale indoor scene point clouds with MASH.

$$\text{Chamfer}(P, Q) = \frac{1}{|P|} \sum_{p \in P} d(p, Q) + \frac{1}{|Q|} \sum_{q \in Q} d(q, P) \quad (5)$$

$$F\text{-Score}(P, Q) = \frac{2 \cdot TP(P, Q) \cdot TP(Q, P)}{TP(P, Q) + TP(Q, P)} \quad (6)$$

$$S_{\text{cos}}(P, Q) = \frac{1}{|P|} \left| \sum_{p \in P} d_n(p, Q) \right| + \frac{1}{|Q|} \left| \sum_{q \in Q} d_n(q, P) \right| \quad (7)$$

where the true-positive function TP is defined as:

$$TP(P, Q) = \frac{1}{|P|} \sum_{p \in P} \mathbb{I}(d(p, Q) \leq 0.01) \quad (8)$$

and the normal similarity function d_n is defined as:

$$d_n(p, Q) = v(p) \cdot v(\arg \min_{q \in Q} (||p - q||_2)) \quad (9)$$

with $v(p)$ to be the surface normal at point p .

When calculating all these metrics, we sample 50k points by farthest point sampling from the ground truth and reconstructed surfaces, respectively. In order to make metrics more intuitive, we

sample the object surface twice with different initial points by and report all metrics in Table 5 as *FPS*. This setting reflects the approximate upper or lower bound of each metric.

Surface reconstruction comparison. For object surface reconstruction, we conduct experiments on the ShapeNet-V2 dataset [Chang et al. 2015] and list detailed comparison results on the seven main categories with the metrics of CD and F-Score in Table 5. Note that since the results of D_H and S_{cos} on different categories are relatively similar, we mainly report their results on the complete dataset in the main paper. More visual examples are also presented in Figure 2.

Surface reconstruction on noisy data. We also conduct experiments on noisy inputs across all categories, in which the noise conforms to a Gaussian distribution with the mean of 0 and the variance of $0.2\%L$ or $0.5\%L$ where L is the length of the bounding box diagonal corresponding to each object. The quantitative comparisons are shown in Table 6, and our method gets consistently better results.

Some representative surface reconstruction results of each method on the Chair category are presented in Figure 3. The global consistency of the normal directions estimated by PGR will decrease as the noise increases, resulting in obvious holes in the results due to the flipping of normal directions. Since ARONet and ConvONet are learning-based methods, they have stronger anti-noise capabilities. However, their methods can still be misled by noise that does not belong to the object surface, leading to incorrect judgments of the object topology. Different from their methods, since the surface of the object is approximately distributed at the mean of the noise data, our method moves the surface patch of each anchor to near the mean position of local noise data by trying to fit all given points and ultimately achieves higher reconstruction quality.

3 SHAPE GENERATION

Technical details. We use the optimal transport conditional flow matching [Tong et al. 2023] to train our diffusion models on 8 RTX4090 with a batch size of 1024 for $T = 1,000$ epochs. The learning rate is linearly increased to $lr_{max} = 2e-4$ in the first $t_0 = 80$ epochs, and then gradually decreases using the cosine decay schedule $lr_{max} * 0.5^{1+\cos(\frac{t-t_0}{T-t_0})}$ until reaching $1e-6$.

To further extract the mesh from the generated MASH representation, we train an extra network to decode the occupancy grid and then extract the mesh by using ODC [Hwang and Sung 2024]. Meanwhile, to enable our shape decoder to adapt to the potentially noisy generated MASH data, we borrow the key idea of VAE [Pinheiro Cinelli et al. 2021]. Specifically, We map each channel of all anchors in the dataset to an approximate normal distribution by carefully designing a strictly monotonically increasing piecewise linear function. This allows us to perform variational MASH sampling more effectively and obtain a more stable shape decoder.

Compared to Shape2VecSet [Zhang et al. 2023], both the training and sampling time of our model is less than one-third, with the same network backbone. More specifically, Shape2VecSet takes about 28 GPU days to train the AutoEncoder and another 92 GPU days to train the diffusion model, while our model takes about 32 GPU days to train with the same setting. When generating shapes on a single RTX 4090 GPU, Shape2VecSet takes about 0.2 seconds to generate

a vector set and about 14 seconds to decode it into a 3D shape. In contrast, our model takes 0.05 seconds to generate a MASH, and another 4 seconds to convert the MASH into a triangular mesh if required.

Evaluation metrics. The metrics used to measure the visual similarity between generated shapes and the dataset are defined as:

$$\text{Render-FID} = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{\frac{1}{2}}) \quad (10)$$

$$\text{Render-KID} = \text{MMD} \left(\frac{1}{|\mathcal{R}|} \sum_{x \in \mathcal{R}} \max_{y \in \mathcal{G}} D(x, y) \right)^2 \quad (11)$$

where g and r denote the generated and training datasets respectively. μ and Σ are the statistical mean and covariance matrix of the feature distribution extracted by the Inception-V3 network [Szegedy et al. 2016]. $D(x, y)$ is a polynomial kernel function to evaluate the similarity of two samples, \mathcal{G} and \mathcal{R} are feature distributions of the generated set and reference set, respectively. The function $\text{MMD}(\cdot)$ is Maximum Mean Discrepancy.

To measure the geometric similarity between the generated 3D shapes and the GT shapes, we sample 50K points from each and use CD, EMD, and F-Score as the primary evaluation metrics. To measure the alignment between images and shapes, we sample 8,192 points from the generated shapes and use ULIP-2 [Xue et al. 2024] as the primary evaluation metric. Specifically, ULIP-I is defined as $\text{ULIP-I}(I, S) = \langle E_I, E_S \rangle$, corresponding to the inner product of normalized ULIP features of image I and generated shape S .

Multi-modal shape generation. In addition to category-conditioned and image-conditioned generative models, we also trained a multi-modal generative model using the ShapeNet-V2 and ShapeGlot [Achlioptas et al. 2019] datasets.

Specifically, we use the ULIP-2 pre-trained model to encode three different modalities including text, images, and point clouds, to obtain the corresponding feature embeddings and serve them as conditions of the generative model. The text label comes from the ShapeGlot dataset, but only for the chair category. For the image condition, we create an image dataset by rendering all shapes in the ShapeNet-V2 dataset from 12 different viewpoints. For the point cloud condition, we sample a point cloud with 8,192 points through the farthest point sampling for each shape in ShapeNet-V2.

The generated results are shown in Figure 4, 5, and 6. Our multi-modal generative model can accurately generate shapes that meet our expectations based on different modalities of input.

Image-conditioned shape generation. We show more visual comparisons of the image-conditioned shape generation results in Figure 7 and more our results in Figure 8.

REFERENCES

- Panos Achlioptas, Judy Fan, Robert Hawkins, Noah Goodman, and Leonidas J Guibas.
2019. ShapeGlot: Learning language for shape differentiation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 8938–8947.
Angel X. Chang, Thomas A. Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qi-Xing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. 2015. ShapeNet: An Information-Rich 3D Model Repository. *CoRR* abs/1512.03012 (2015), arXiv:1512.03012 <http://arxiv.org/abs/1512.03012>
Philipp Erler, Paul Guerrero, Stefan Ohrhallinger, Niloy J. Mitra, and Michael Wimmer. 2020. Points2Surf: Learning Implicit Surfaces from Point Clouds. In *Computer Vision*

Table 5. Detailed quantitative comparison with surface reconstruction baselines, including SPR+PCA [Kazhdan and Hoppe 2013], PGR [Lin et al. 2022], ARO-Net [Wang et al. 2023] and ConvONet [Peng et al. 2020] on the ShapeNet dataset. For ease of comparison of results, we multiply the L1 Chamfer Distance by 1000.

		SPR+PCA	PGR	ConvONet	ARONet	Ours	FPS
Chamfer ↓	airplane	63.562	4.839	13.254	11.569	4.194	3.426
	chair	102.852	7.824	19.374	15.713	6.646	5.471
	car	135.148	8.052	25.073	23.173	6.795	5.383
	lamp	88.064	6.598	16.903	17.958	4.534	3.850
	rifle	79.631	4.479	11.393	13.171	3.829	2.916
	sofa	112.850	7.857	26.782	18.069	6.444	5.249
	table	77.971	7.595	22.339	15.965	6.578	5.473
		mean	89.565	6.381	17.732	15.697	5.450
F-Score ↑	airplane	0.544	0.997	0.871	0.873	0.998	1.000
	chair	0.436	0.984	0.801	0.886	0.997	1.000
	car	0.391	0.964	0.627	0.720	0.987	0.998
	lamp	0.517	0.979	0.811	0.801	0.999	0.999
	rifle	0.525	0.987	0.931	0.924	0.998	1.000
	sofa	0.422	0.970	0.698	0.727	0.996	0.999
	table	0.530	0.988	0.802	0.892	0.997	0.999
		mean	0.497	0.988	0.812	0.880	0.997
$D_H \downarrow$		mean	0.272	0.023	0.131	0.117	0.019
$S_{cos} \uparrow$		mean	0.684	0.974	0.821	0.898	0.980

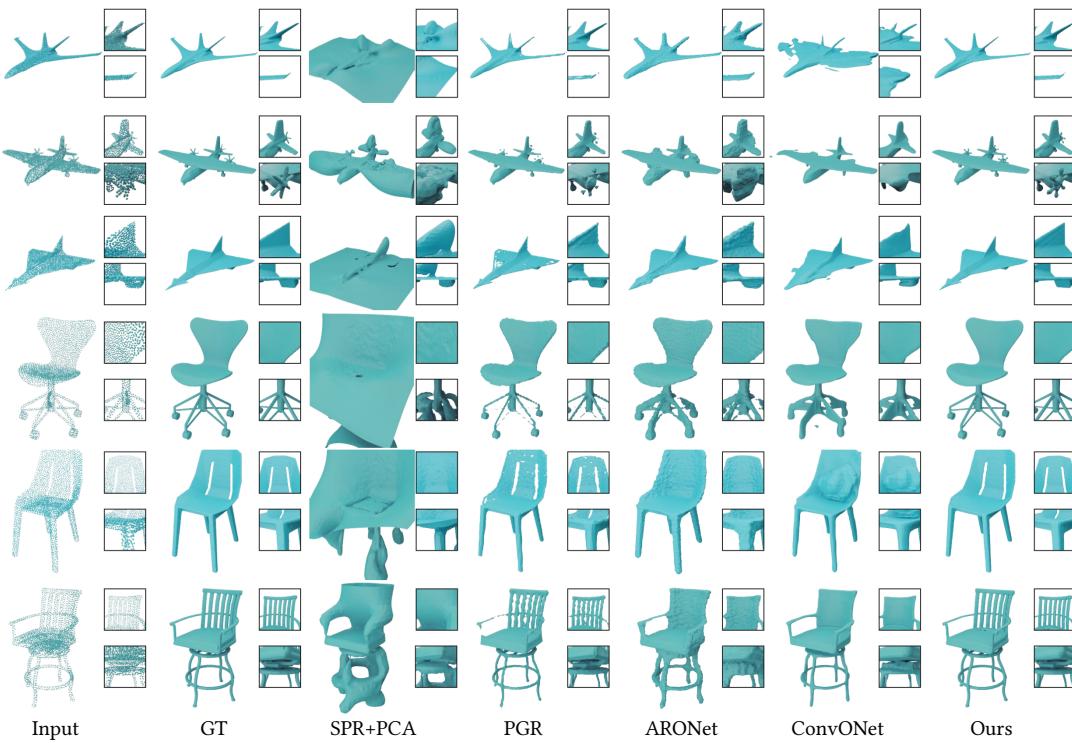


Fig. 2. Qualitative results on surface reconstruction with different methods. We additionally select two details for each result to show the performance of all methods better.

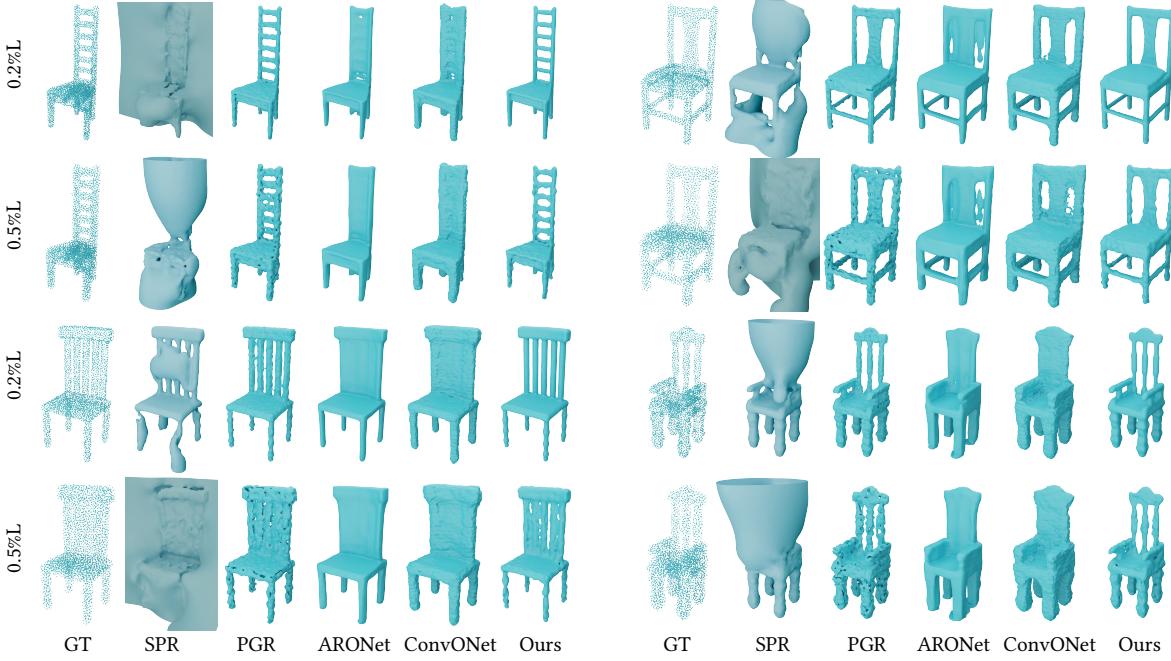


Fig. 3. Reconstructions on the chair category with two different noise levels.

Table 6. Quantitative comparison with surface reconstruction baselines on noisy input. For ease of comparison of results, we multiply the L1 Chamfer Distance by 1000 here.

	SPR+PCA	PGR	ConvONet	ARONet	Ours
0.2%L	Chamfer ↓	89.724	7.752	18.203	16.825
	F-Score ↑	0.494	0.949	0.803	0.878
	$D_H \downarrow$	0.271	0.024	0.142	0.122
	$S_{cos} \uparrow$	0.542	0.966	0.803	0.884
0.5%L	Chamfer ↓	94.172	19.875	21.434	18.872
	F-Score ↑	0.493	0.802	0.781	0.804
	$D_H \downarrow$	0.268	0.069	0.148	0.125
	$S_{cos} \uparrow$	0.326	0.828	0.801	0.881

- ECCV 2020, Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (Eds.). Springer International Publishing, Cham, 108–124.
- H. Fan, H. Su, and L. Guibas. 2017. A Point Set Generation Network for 3D Object Reconstruction from a Single Image. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, Los Alamitos, CA, USA, 2463–2471. <https://doi.org/10.1109/CVPR.2017.264>
- D. P. Huttenlocher, G. A. Klanderman, and W. A. Rucklidge. 1993. Comparing Images Using the Hausdorff Distance. *IEEE Trans. Pattern Anal. Mach. Intell.* 15, 9 (sep 1993), 850–863. <https://doi.org/10.1109/34.232073>
- Jisung Hwang and Minhyuk Sung. 2024. Occupancy-Based Dual Contouring. In *SIGGRAPH Asia 2024 Conference Papers*. 1–11.
- Michael Kazhdan and Hugues Hoppe. 2013. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)* 32, 3 (2013), 1–13.
- Siyou Lin, Dong Xiao, Zuochang Shi, and Bin Wang. 2022. Surface Reconstruction from Point Clouds without Normals by Parametrizing the Gauss Formula. *ACM Trans. Graph.* 42, 2, Article 14 (oct 2022), 19 pages. <https://doi.org/10.1145/3554730>
- I Loshchilov. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).
- Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. 2019. DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation. *arXiv:1901.05103 [cs.CV]*
- Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. 2020. Convolutional Occupancy Networks. In *Computer Vision – ECCV 2020*, Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (Eds.). Springer International Publishing, Cham, 523–540.

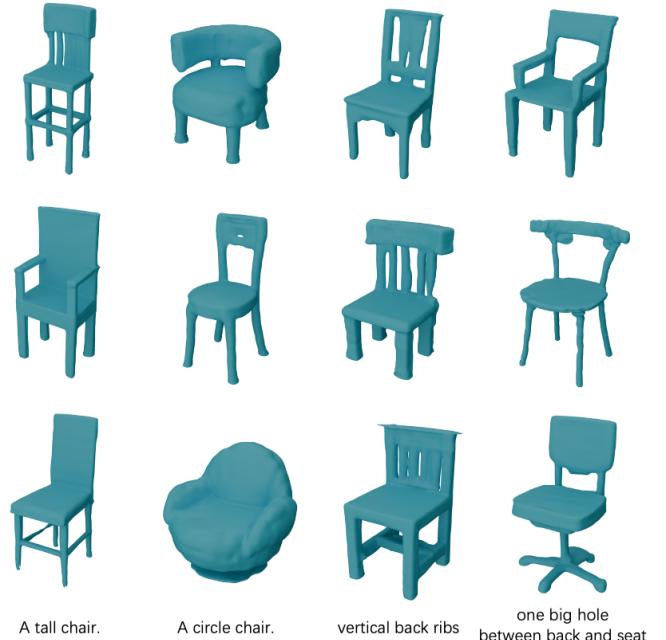


Fig. 4. Qualitative results on text-conditioned generation.



Towards scalable multimodal pre-training for 3d understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 27091–27101. Biao Zhang, Jiapeng Tang, Matthias Niessner, and Peter Wonka. 2023. 3dshape2vecset: A 3d shape representation for neural fields and generative diffusion models. *ACM Transactions on Graphics (TOG)* 42, 4 (2023), 1–16.

Fig. 5. Qualitative results on image-conditioned generation.

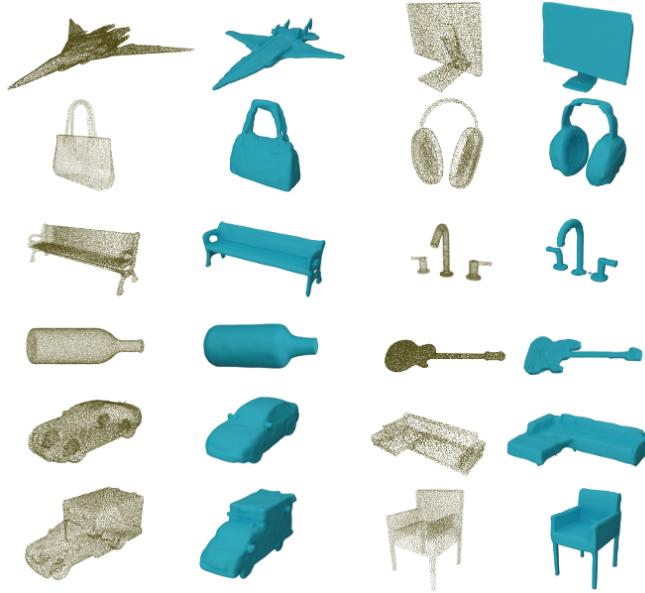


Fig. 6. Qualitative results on point cloud-conditioned generation.

- Lucas Pinheiro Cinelli, Matheus Araújo Marins, Eduardo Antônio Barros da Silva, and Sérgio Lima Netto. 2021. Variational autoencoder. In *Variational Methods for Machine Learning with Applications to Deep Networks*. Springer 111–149.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2818–2826.
- Alexander Tong, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrid Rector-Brooks, Kilian Fatras, Guy Wolf, and Yoshua Bengio. 2023. Conditional flow matching: Simulation-free dynamic optimal transport. *arXiv preprint arXiv:2302.00482* 2, 3 (2023).
- Yizhi Wang, Zeyu Huang, Ariel Shamir, Hui Huang, Hao Zhang, and Ruizhen Hu. 2023. ARO-Net: Learning Implicit Fields from Anchored Radial Observations. *CVPR* (2023).
- Le Xue, Ning Yu, Shu Zhang, Artemis Panagopoulou, Junnan Li, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, et al. 2024. Ulip-2:

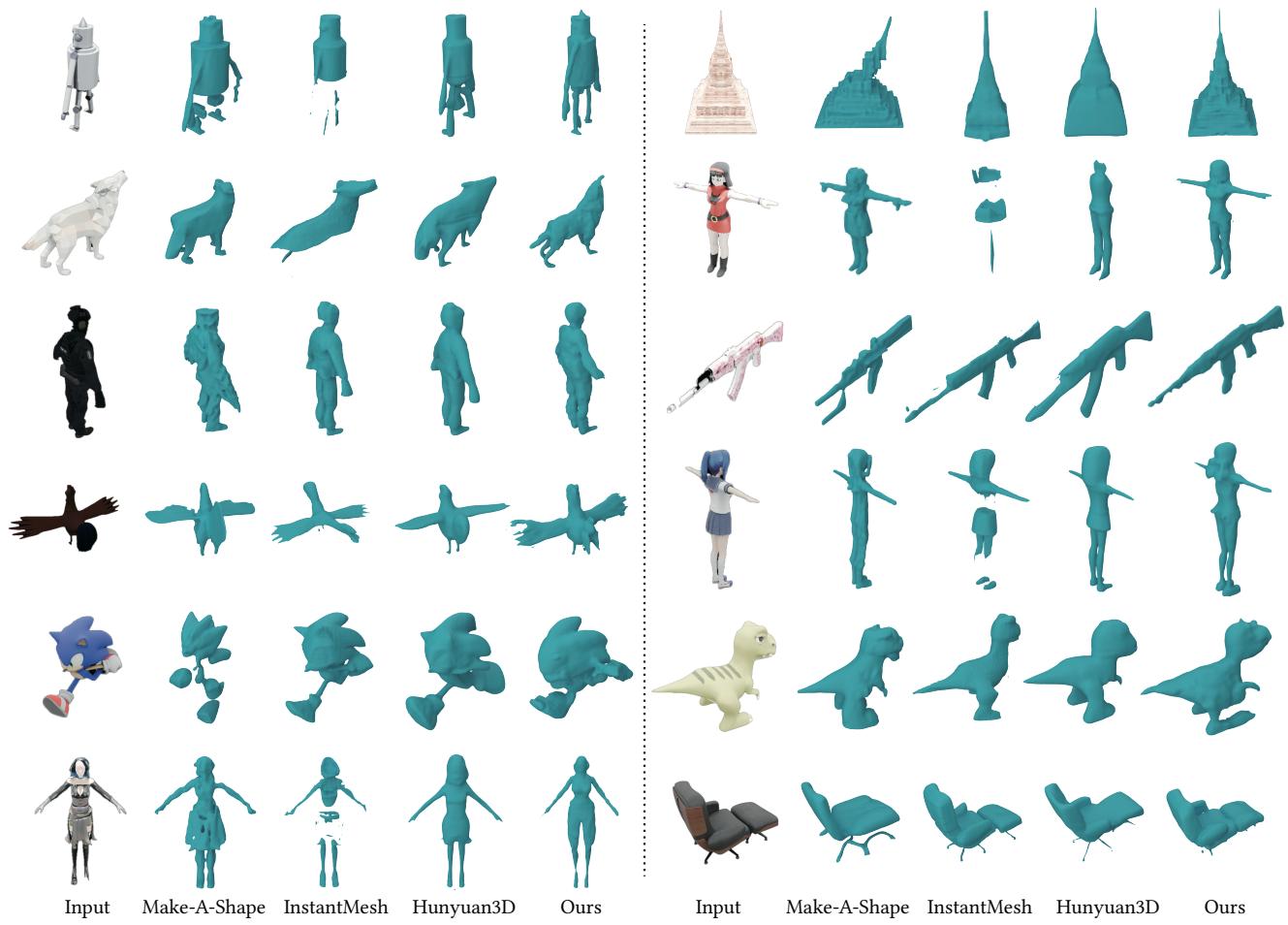


Fig. 7. More qualitative results on single image to 3D generation compared with different methods.



Fig. 8. More qualitative results on single image to 3D generation based on MASH.