

MeshAnything: Artist-Created Mesh Generation with Autoregressive Transformers

Yiwen Chen^{1,2*}, Tong He^{2†}, Di Huang², Weicai Ye², Sijin Chen³, Jiayang Tang⁴
Xin Chen⁵, Zhongang Cai⁶, Lei Yang⁶, Gang Yu⁷, Guosheng Lin^{1†}, Chi Zhang^{8†}

¹S-Lab, Nanyang Technological University ²Shanghai AI Lab

³Fudan University ⁴Peking University ⁵University of Chinese Academy of Sciences

⁶SenseTime Research ⁷Stepfun ⁸Westlake University

<https://buaacyw.github.io/mesh-anything/>

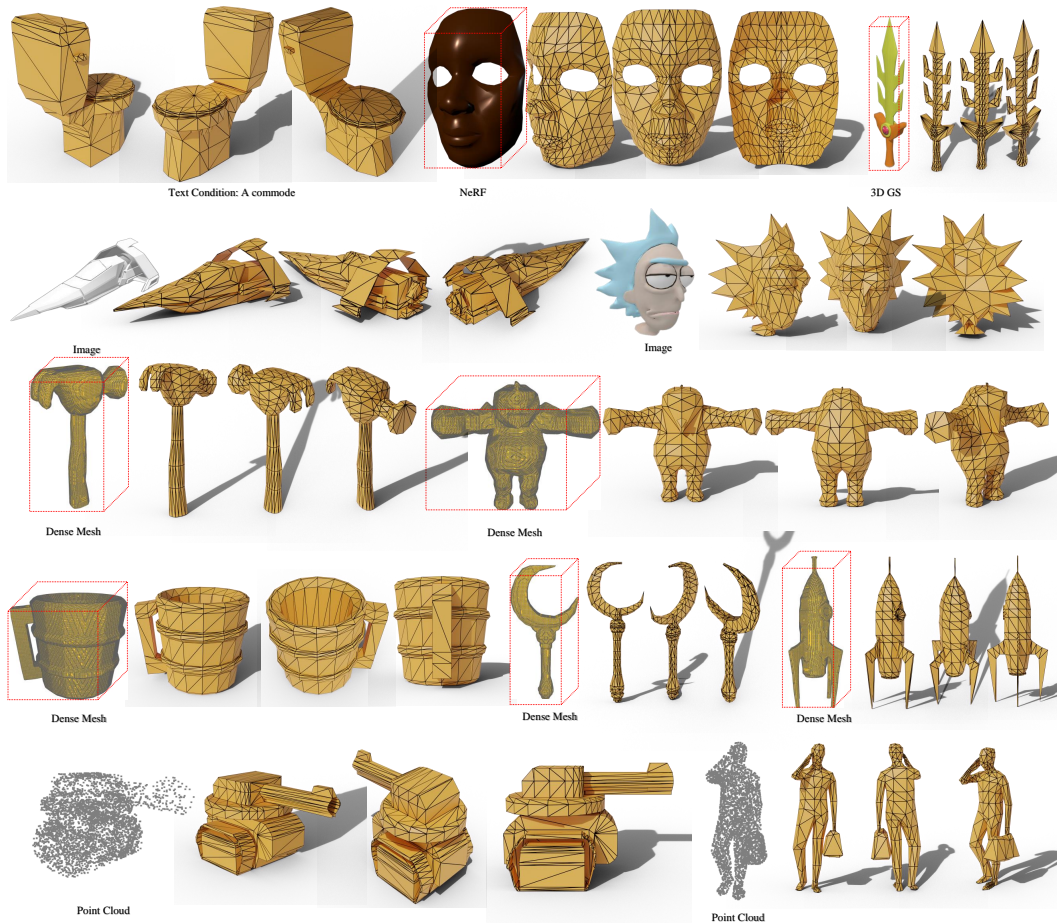


Figure 1: **MeshAnything converts any 3D representation into Artist-Created Meshes (AMs), i.e., meshes created by human artists.** It can be combined with various 3D asset production pipelines, such as 3D reconstruction and generation, to transform their results into AMs that can be seamlessly applied in the 3D industry.

*Work done during a research internship at Shanghai AI Lab.

†Corresponding Authors.

Abstract

Recently, 3D assets created via reconstruction and generation have matched the quality of manually crafted assets, highlighting their potential for replacement. However, this potential is largely unrealized because these assets always need to be converted to meshes for 3D industry applications, and the meshes produced by current mesh extraction methods are significantly inferior to Artist-Created Meshes (AMs), i.e., meshes created by human artists. Specifically, current mesh extraction methods rely on dense faces and ignore geometric features, leading to inefficiencies, complicated post-processing, and lower representation quality. To address these issues, we introduce MeshAnything, a model that treats mesh extraction as a generation problem, producing AMs aligned with specified shapes. By converting 3D assets in any 3D representation into AMs, MeshAnything can be integrated with various 3D asset production methods, thereby enhancing their application across the 3D industry. The architecture of MeshAnything comprises a VQ-VAE and a shape-conditioned decoder-only transformer. We first learn a mesh vocabulary using the VQ-VAE, then train the shape-conditioned decoder-only transformer on this vocabulary for shape-conditioned autoregressive mesh generation. Our extensive experiments show that our method generates AMs with hundreds of times fewer faces, significantly improving storage, rendering, and simulation efficiencies, while achieving precision comparable to previous methods.

1 Introduction

In recent years, the 3D community has experienced rapid advancements, with a variety of methods developed for automatically producing high-quality 3D assets. These methods, including 3D reconstruction [44, 72, 2, 3, 34, 31], 3D generation [47, 39, 64, 41, 56, 30, 57, 69, 66], and scanning [17, 26, 27], can produce 3D assets with shape and color quality comparable to manually created ones. The success of these methods reveals the potential to replace manually created 3D models with automatically produced ones in the 3D industry, including applications in games, movies, and the metaverse, significantly reducing time and labor costs.

However, this potential remains largely unrealized because the current 3D industry predominantly relies on mesh-based pipelines for their superior efficiency and controllability, while methods for producing 3D assets typically use alternative 3D representations to achieve optimal results across various scenarios. Therefore, substantial efforts [42, 16, 43, 51, 12, 53] are devoted to converting other 3D representations into meshes and have achieved some success. Meshes produced by these methods approximate the shape quality of those created by human artists, which we refer to as Artist-Created Meshes (AMs), but they still fall short in addressing the aforementioned issues.

This is because all meshes produced by these methods [42, 16, 43, 51, 12, 53] exhibit significantly poorer topology quality compared to AMs. As shown in Fig. 2, these methods rely on dense faces to reconstruct 3D shapes, completely ignoring geometric characteristics. Using these meshes in the 3D industry leads to three significant problems: First, converted meshes typically contain several orders of magnitude more faces compared to AMs, leading to significant inefficiencies in storage, rendering, and simulation. Moreover, the converted meshes complicate post-processing and downstream tasks in the 3D pipeline. They significantly increase the challenge for human artists in optimizing these meshes due to their chaotic and inefficient topologies. Finally, previous methods struggle to represent sharp edges and flat surfaces, resulting in oversmoothing and bumpy artifacts as shown in Fig. 2.

In this work, we aim to solve the aforementioned issues to facilitate the application of automatically generated 3D assets in the 3D industry. As mentioned earlier, all previous methods [42, 16, 43, 51, 12, 53] extract 3D meshes with excessively dense faces in a reconstruction manner, which inherently cannot solve these issues. Therefore, we diverge from previous approaches by formulating mesh extraction as a generation problem for the first time: we teach models to generate Artist-Created Meshes (AMs) that are aligned with the given 3D assets. The meshes generated by our methods mimic the shape and topology quality of those created by human artists. Consequently, our setting, namely Shape-Conditioned AM Generation, is fundamentally free from all previous issues, enabling seamless integration of the generated results into the 3D industry pipeline.

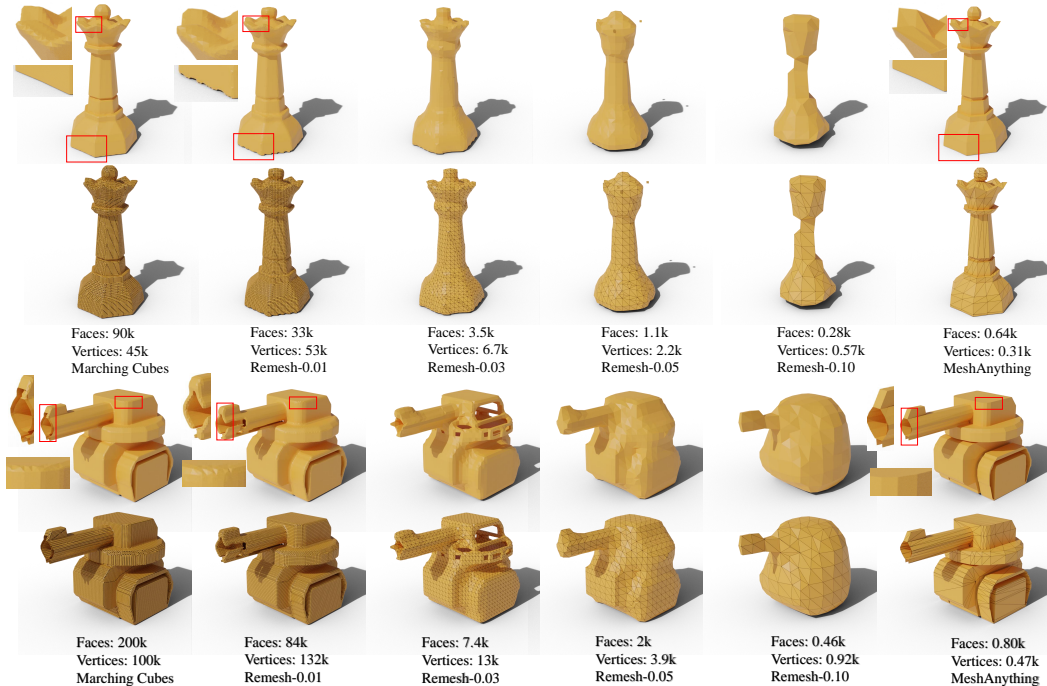


Figure 2: **Comparison with Marching Cubes [42] and Remesh [4].** We apply Marching Cubes and MeshAnything to ground truth shapes and then apply remeshing to the Marching Cubes results with different voxel sizes. Existing methods extract meshes in a reconstruction manner, ignoring the geometric features of the object and producing dense meshes with poor topology. These methods fundamentally fail to capture sharp edges and flat surfaces, as shown in the zoomed-in figure.

However, training such a model presents significant challenges. The first challenge is constructing the dataset, as we need paired shape conditions and Artist-Created Meshes (AMs) for model training. The shape condition must be efficiently derived from as many diverse 3D representations as possible to serve as a condition during inference. Additionally, it must have sufficient precision to accurately represent 3D shapes and be efficiently processed into features that can be injected into the model. After weighing the trade-offs, we chose point clouds due to their explicit and continuous representation, ease of derivation from most 3D representations, and the availability of mature point cloud encoders [48, 49, 75].

We filter out high-quality AMs from Objaverse [19, 18] and ShapeNet [7]. When obtaining paired shape conditions, a naive approach would be to sample point clouds directly from AMs. However, this leads to poor results during inference because the sampled point clouds have excessive precision, while automatically produced 3D assets cannot provide point clouds of similar quality, causing a domain gap between training and inference. To address this issue, we intentionally corrupt the shape quality of AMs. We first extract the signed distance function from AMs [63], convert it into a coarser mesh using [42], and then sample point clouds from this coarse mesh to narrow the domain gap in shape conditions between inference and training.

Following [55], we use a VQ-VAE [61] to learn a mesh vocabulary and train a decoder-only transformer [62] on this vocabulary for mesh generation. To inject shape condition, we draw inspiration from the recent success of multimodal large language models (MLLM) [68, 37], where image features encoded by pre-trained image encoders are projected into the token space of the large language models for efficient multimodal understanding. Similarly, we treat the mesh tokens obtained from the trained VQ-VAE as the language token in LLMs and use a pre-trained encoder [75] to encode the point clouds into shape features, which is later projected into the mesh token space. These shape tokens are placed at the beginning of the mesh token sequences, effectively serving as the shape conditions for next-token predictions. After predictions, these predicted mesh tokens are decoded back to meshes with the VQ-VAE decoder [55].

To further enhance the quality of mesh generation, we develop a novel noise-resistant decoder for robust mesh decoding. Our observation is that as the decoder in the VQ-VAE [61] is only trained with ground truth token sequences from the encoder, it could potentially lead to a domain gap when decoding the generated token sequences. To mitigate this problem, we inject the shape condition into the VQ-VAE decoder as auxiliary information for robust decoding and fine-tune it after the VQ-VAE training. This fine-tuning process involves adding noise to the mesh token sequences to simulate possible poor-quality token sequences from the decoder-only transformer, thus making the decoder robust to such poor-quality sequences.

Finally, we introduce our model, MeshAnything, trained based on the aforementioned techniques. As shown in Fig. 1, MeshAnything can convert 3D assets across various 3D representations into AMs, thereby significantly facilitating their application. Furthermore, our extensive experiments demonstrate that our method generates AMs with significantly fewer faces and more refined topology, while achieving precision metrics that are close to or comparable with previous methods.

In summary, our contributions are threefold:

- We highlight one important reason why current automatically produced 3D assets cannot replace those created by human artists: current methods cannot convert these 3D assets into Artist-Created Meshes (AMs). To solve this issue, we propose a novel solution called Shape-Conditioned AM Generation, which aims to generate AMs aligned with given shapes.
- We introduce MeshAnything for Shape-Conditioned AM Generation. MeshAnything can be integrated with various 3D asset production methods, converting their results into AMs to facilitate their application in the 3D industry.
- We develop a novel noise-resistant decoder to enhance mesh generation quality. We inject the shape condition into the decoder as auxiliary information for robust decoding and fine-tune it using noised token sequences to narrow the domain gap between training and inference.

2 Related Works

2.1 Mesh Extraction

Methods for extracting meshes from 3D models are numerous and have been a subject of research for decades. Following [53], we categorize these methods into two main types: Isosurface Extraction [42, 5, 16, 6, 43, 12] and Gradient-Based Mesh Optimization [9, 23, 28, 33, 52, 36, 53].

Traditional isosurface extraction methods [42, 43, 16, 21, 32, 50, 13, 12] focus on extracting a polygonal mesh that represents the level set of a scalar function, an area that has seen extensive study in various fields. The most popular method among them is Marching Cubes [42]. It divides the space into cells, within which polygons are created to approximate the surface. Marching Cubes has been widely used for mesh extraction its robustness and simplicity. To more effectively capture sharp features, Dual Contouring [50] adopts a dual representation by extracting mesh vertices per cell and estimates vertex positions based on local isosurface details. Dual Marching Cubes [46] is another advanced approach that combines the advantages of both Marching Cubes and Dual Contouring. Recently, [13] and [12] introduce a data-driven method to determine the position of the extracted mesh based on the input field.

Transitioning to more recent developments, the advent of machine learning has ushered in new techniques for generating 3D meshes [9, 23, 28, 33, 52, 36, 53]. This line of work explores using neural networks to generate 3D meshes, where the network parameters are optimized through gradient-based methods under specific loss functions. These approaches leverage the computational power of machine learning to refine mesh quality and adaptability, making them suitable for more complex geometrical structures that traditional methods may struggle to process efficiently. [23] proposes predicting a deformable tetrahedral grid for 3D shape representation. Besides, [52] employs a differentiable Marching Tetrahedra layer for mesh extraction. Similar to [52], [53] iteratively optimizes a 3D surface mesh by representing it as the isosurface of a scalar field.

However, these approaches fundamentally differ from ours. They ignore the characteristics of the shape and inherently cannot produce meshes with efficient topology. In contrast, Shape-Conditioned AM Generation formulates mesh extraction as a generation problem for the first time, aiming to

mimic human artists in mesh extraction and thereby generating Artist-Created Meshes (AMs) with hundreds of times fewer faces.

2.2 3D Mesh Generations

3D mesh generation can be mainly divided into two categories: generating dense meshes similar to those produced by previous mesh extraction methods, and generating Artist-Created Meshes (AMs).

The former category is currently the mainstream research focus. Methods such as [24, 66, 69] directly generate meshes in a feed-forward manner, but because they produce dense meshes with low-quality topology similar to previous mesh extraction methods, they still encounter the same issues when applied in the 3D industry. The most popular techniques in this category are LRM-based mesh generation models [30, 66, 69], which utilize transformers to generate 3D meshes in a feed-forward manner, often conditioning on images for control.

Notably, numerous 3D generation methods [47, 59, 64, 11, 58, 71, 30, 22, 10, 38, 54, 35, 14, 15, 57, 65, 60] can also produce meshes. These methods do not directly generate meshes but instead produce 3D assets in other representations. In practice, these 3D assets often need to be converted into meshes using mesh extraction methods, establishing a connection to our approach. Similar to direct mesh generation methods, these approaches primarily focus on the generation of shape and color while ignoring the topology of the mesh, resulting in dense meshes. Consequently, they face challenges when applied to the 3D industry due to their inefficient topology.

Recently, several works have focused on the second category: generating Artist-Created Meshes (AMs) [45, 1, 55, 8]. Although our approach also focuses on AM generation, it fundamentally differs from these methods. Since they lack shape conditioning, these methods must simultaneously learn the complex 3D shape distribution—which typically alone requires extensive training [30, 57]—and the topology distribution of AMs, leading to very challenging training processes. In contrast, Shape-Conditioned AM Generation only needs to learn how to construct efficient topology for a known shape, making the learning process significantly easier. Numerous mature methods [44, 34, 30, 57, 2, 3, 66, 69] already effectively produce high-quality 3D shapes; by combining our approach with these methods, we can achieve similar outcomes to those of direct AM generation methods. This will be detailed in Sec. 3.

Among these methods, the most relevant to ours is MeshGPT [55], as we follow its architecture. [55] introduced a combination of a VQ-VAE [61] and an autoregressive transformer architecture. It first learns a mesh vocabulary with the VQ-VAE and then trains the transformer on the learned vocabulary for mesh generation. However, MeshGPT’s results are limited to several categories in ShapeNet. MeshGPT requires a training GPU hours similar to ours, but our method can generalize to unlimited categories in Objaverse. This is largely due to the difference in target complexity caused by MeshGPT needing to additionally learn the complex 3D shape distribution.

3 Shape-Conditioned AM Generation

In this section, we first introduce the formal formulation for Shape-Conditioned AM Generation and compare it with previous mesh generation settings [45, 55, 1]. We show that it can achieve the same targets as previous mesh generation methods with significantly less training effort.

Shape-Conditioned AM Generation targets to estimate a conditional distribution $p(\mathcal{M}|\mathcal{S})$. In this formula, \mathcal{M} refers to the Artist-Created Mesh (AM), i.e., the mesh manually modeled by human artists. \mathcal{S} refers to the 3D shape information that indicates the 3D shape to which \mathcal{M} should align. The input form of \mathcal{S} can be diverse, such as voxels or point clouds. Therefore, this versatility allows our method to be integrated with any 3D pipeline that outputs \mathcal{S} , such as 3D reconstruction [44, 34], generation [47, 30], and scanning, making these methods more efficient for the 3D industry.

Compared to existing AM generation work, they directly estimate the distribution $p(\mathcal{M}|\mathcal{C})$, where \mathcal{C} denotes conditions such as images, text or empty sets for unconditional generation. However, estimating $p(\mathcal{M}|\mathcal{C})$ requires an understanding of both the underlying shape, i.e., \mathcal{S} , and complex topological structures \mathcal{M} . Given this, we made the following approximation:

$$p(\mathcal{M}|\mathcal{C}) \approx p(\mathcal{M}, \mathcal{S}|\mathcal{C}). \tag{1}$$

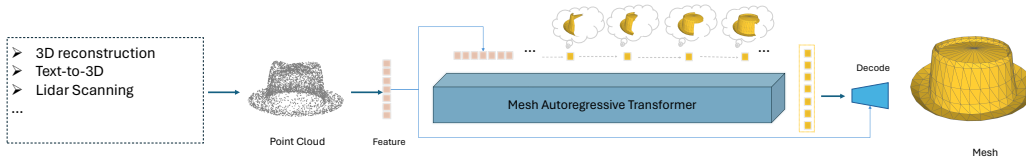


Figure 3: **Pipeline Overview.** We introduce MeshAnything, an autoregressive transformer capable of generating Artist-Created Meshes that adhere to given 3D shapes. We sample point clouds from given 3D assets, encode them into features, and inject them into the decoder-only transformer to achieve shape-conditional mesh generation.

According to the chain rule, we have:

$$p(\mathcal{M}, \mathcal{S}|\mathcal{C}) = p(\mathcal{M}|\mathcal{S}, \mathcal{C}) \cdot p(\mathcal{S}|\mathcal{C}). \quad (2)$$

For distribution $p(\mathcal{M}|\mathcal{S}, \mathcal{C})$, given that \mathcal{S} is a much stronger and more direct condition than \mathcal{C} , we can make the following approximation:

$$p(\mathcal{M}|\mathcal{S}, \mathcal{C}) \approx p(\mathcal{M}|\mathcal{S}). \quad (3)$$

Combining 1, 2 and 3:

$$p(\mathcal{M}|\mathcal{C}) \approx p(\mathcal{M}|\mathcal{S}) \cdot p(\mathcal{S}|\mathcal{C}), \quad (4)$$

in which $p(\mathcal{M}|\mathcal{S})$ is the focus of our shape-conditioned mesh generation. In the 3D community, numerous large models [30, 57, 69, 55] aim to estimate $p(\mathcal{S}|\mathcal{C})$ using various 3D representations and demonstrate excellent results. Besides, some single scene 3D asset production methods [44, 34, 2, 3, 47, 40, 56] can also provide samples from this distribution. By integrating our framework with these existing methods, we can leverage their capabilities to enhance our mesh generation process. This integration allows for a more resource-efficient way to estimate $p(\mathcal{M}|\mathcal{C})$, significantly reducing the complexity and resources required compared to previous methods.

4 Method

In this section, we detail our shape condition strategy in Section 4.1. After that, we provide a detailed description for MeshAnything, which consists of a VQVAE with our newly proposed noise-resistant decoder (Section 4.2) and a shape-conditioned autoregressive transformer (Section 4.3).

4.1 Shape Encoding for Conditional Generation

We begin by describing our shape condition strategy. MeshAnything targets learning $p(\mathcal{M}|\mathcal{S})$, so we need to pair each mesh \mathcal{M} with a corresponding \mathcal{S} , i.e., the shape condition. Choosing an appropriate 3D representation for \mathcal{S} is non-trivial and should satisfy the following conditions:

1. It should be easily extracted from various 3D representations. This ensures that the trained models can be integrated with a wide range of 3D asset production pipelines [44, 34, 30, 47, 57].
2. It should be suitable for data augmentation to prevent overfitting. To ensure the effectiveness of \mathcal{S} during training, any data augmentation applied to \mathcal{M} must be equivalently applicable to \mathcal{S} .
3. It should be efficiently and conveniently input into the model as a condition. To ensure the model comprehends the shape information and to maintain efficient training, \mathcal{S} must be easily and effectively encoded into features.

Considering the first and second points, \mathcal{S} should be in an explicit representation. Further considering the third point, the main explicit 3D representations that can be easily encoded as features are voxels and point clouds. Both representations are suitable, but voxels typically require a high resolution to accurately represent shapes, and processing high-resolution voxels into features is computationally expensive. Additionally, voxels, being a discrete representation, are less precise for data augmentation

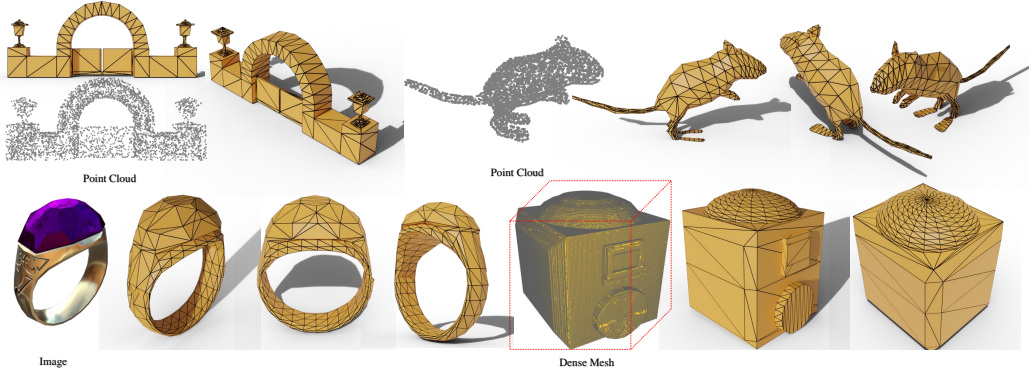


Figure 4: **Additional qualitative results of MeshAnything.** As shown, MeshAnything can be integrated with various 3D production pipelines to achieve highly controllable mesh generation.

compared to point clouds. Therefore, we chose point clouds as the representation for \mathcal{S} . To enhance the expressive power of the point clouds, we also include normals into the point cloud representation.

To obtain point clouds from the ground truth mesh for training, we could simply sample point clouds directly from the surface of \mathcal{M} . However, this would create problems during inference: the surfaces of automatically generated 3D assets are often rougher than those of AMs. For example, in AMs, we would sample a series of points on a flat plane, whereas automatically generated 3D assets would have uneven surfaces, causing a domain gap between training and inference.

Therefore, we need to ensure that \mathcal{S} extracted from the ground truth \mathcal{M} during training has a similar domain to the \mathcal{S} extracted during inference. To bring their domains closer, we intentionally construct coarse meshes from AMs. We first extract the signed distance function from \mathcal{M} with [63], then convert it into a relatively coarse mesh using Marching Cubes [42] to destroy the ground truth topology. Finally, we sample point cloud and its normals from the coarse mesh. This approach also helps to avoid overfitting, as AMs typically have fewer faces, and each face can often sample multiple points. The network can easily recognize the ground truth topology by determining whether the points lie on the same plane.

Since almost all 3D representations can be converted into a coarse mesh using Marching Cubes [42] or sampled into point clouds, this ensures that the domain of \mathcal{S} is consistent during both training and inference. We pair the point clouds extracted as \mathcal{S} with \mathcal{M} to create a data item $\{(\mathcal{M}_i, \mathcal{S}_i)\}_i$ for training.

4.2 VQ-VAE with Noise-Resistant Decoder

Following MeshGPT [55], we first train a VQ-VAE [61] to learn a vocabulary of geometric embeddings for better transformer [62] learning. Different to MeshGPT, which uses graph convolutional networks [67] and ResNet [29] as the encoder and decoder respectively, we employ transformers with identical structures for both the encoder and decoder. When training VQ-VAE, meshes are discretized and input as a sequence of triangle faces:

$$\mathcal{M} := (f_1, f_2, f_3, \dots, f_N), \quad (5)$$

where f_i is the coordinates of the vertices of each face, and N is the number of faces in \mathcal{M} . The encoder E then extracts a feature vector for each face:

$$\mathcal{Z} = (z_1, z_2, \dots, z_N) = E(\mathcal{M}), \quad (6)$$

where z_i is the feature vector for f_i .

The extracted faces are then quantized into quantized features \mathcal{T} with codebook \mathcal{B} :

$$\mathcal{T} = RQ(\mathcal{Z}; \mathcal{B}) \quad (7)$$

Finally, the reconstructed mesh is decoded from \mathcal{T} with decoder D by predicting the logits for each vertex's coordinates:

$$\hat{\mathcal{M}} = D(\mathcal{Z}) \quad (8)$$

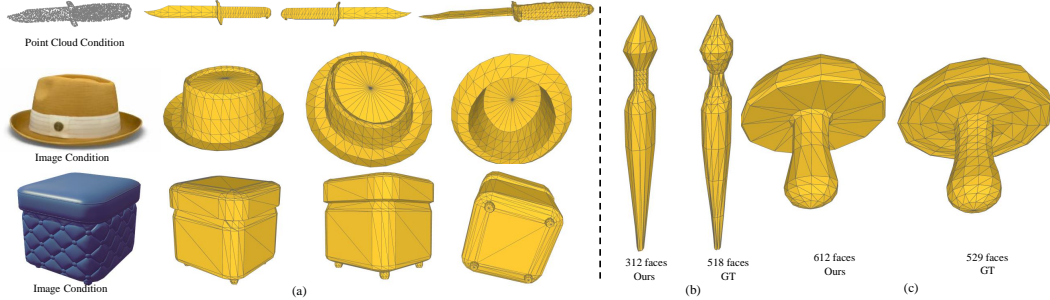


Figure 5: **Qualitative Results.** (a) further demonstrates our capability to achieve highly controllable mesh generation when combined with 3D asset production pipelines. Besides, we compare our results with ground truth in (b) and (c). In (b), MeshAnything generates meshes with better topology and fewer faces than the ground truth. In (c), we produce meshes with a completely different topology while achieving a similar shape, proving that our method does not simply overfit but understands how to construct meshes using efficient topology.

The VQ-VAE is trained end-to-end with cross-entropy loss on the predicted vertex coordinate logits and the commitment loss of vector quantization [61]. After the training of VQ-VAE, the encoder-decoder of VQ-VAE is treated as a tokenizer and detokenizer for autoregressive transformer training.

However, as shown in Fig. 6, there are possible imperfections in the generation results. To address this issue, given our setting of Shape-Conditioned AM Generation, the VQ-VAE decoder can also take the shape condition as input. Small imperfections in the token sequences generated by the transformer can potentially be corrected by a shape-aware decoder. Therefore, after completing the vanilla VQ-VAE training, we add an additional decoder fine-tuning stage, where we inject the shape information into the transformer decoder. Then we add random Gumbel noise to the codebook sampling logits to simulate the potential imperfections in the token sequences generated by the transformer during inference. The decoder is then updated independently with the same cross-entropy loss to train it to produce refined meshes even when facing imperfect token sequences. Our experiments in Tab. 1 and Tab. 2 show that our method effectively enhances the decoder’s noise resistance and mesh generation quality.

4.3 Shape-Conditioned Autoregressive Transformer

To add shape condition to the transformer, inspired by the success of multimodal large language models [68, 37, 70, 25], we first encode the point cloud into a fixed-length token sequence with a point cloud encoder \mathcal{P} and then concatenate it to the front of the embedding sequence from \mathcal{T} VQ-VAE as the final input embedding sequence for the transformer:

$$\mathcal{T}' = \text{concat}(\mathcal{P}(\mathcal{S}), \mathcal{T}) \quad (9)$$

where \mathcal{T}' is the training input for the transformer.

We borrow a pretrained point encoder from [75] and add a linear projection layer to project its output feature to the same latent space as \mathcal{T} . During training, the original point encoder from [75] is frozen; we only update the newly added projection layer and the autoregressive transformer with cross-entropy loss.

During inference, we input $\mathcal{P}(\mathcal{S})$ to the transformer and require it to generate the subsequent sequence, $\hat{\mathcal{T}}$. $\hat{\mathcal{T}}$ is then input to the noise-resistant decoder to reconstruct meshes:

$$\hat{\mathcal{M}} = D(\hat{\mathcal{T}}) \quad (10)$$

where $\hat{\mathcal{M}}$ is the final generated AM.

We use the standard next-token prediction loss to train shape-conditioned transformers. For each sequence, we add a $\langle \text{bos} \rangle$ token after the point cloud tokens and a $\langle \text{eos} \rangle$ token after the mesh tokens to identify the end of a 3D mesh.

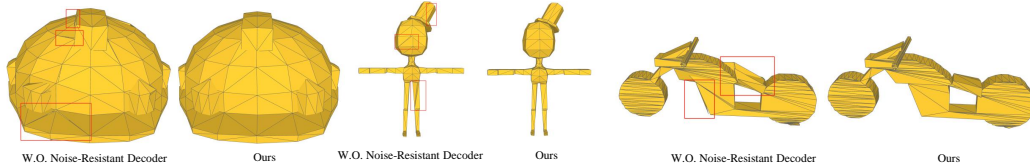


Figure 6: **Ablation on Noise-Resistant Decoder.** The decoder-only transformer may generate low-quality token sequences, and the decoder of VQ-VAE would typically produce flawed meshes based on these sequences. In contrast, our Noise-Resistant Decoder, aided by shape conditions, has the ability to resist these low-quality token sequences, producing higher-quality meshes.

Table 1: Reconstruction Performance under Different Noise Levels with and without Noise-Resistant (NR) Decoder. Please refer to 5.4 for metrics explanation.

Noise Level	$CD(\times 10^{-2})\downarrow$		$ECD(\times 10^{-2})\downarrow$		$NC\uparrow$	
	W/O NR	W/ NR	W/O NR	W/ NR	W/O NR	W/ NR
0.0	0.011	0.007	0.035	0.023	0.987	0.993
0.1	0.187	0.028	0.613	0.138	0.973	0.991
0.5	1.167	0.639	2.538	1.329	0.964	0.981
1.0	2.131	1.798	4.317	2.316	0.952	0.969

5 Experiments

5.1 Data Preparation

Data Selection. Existing AM generation works are limited to a few categories. However, our method targets to operate on general shapes. MeshAnything is trained on a combined dataset of Objaverse [19] and ShapeNet [7], selected for their complementary characteristics. We chose Objaverse because it contains a large number of AMs without category limitations. On the other hand, ShapeNet offers higher data quality within limited categories.

We filter out meshes with more than 800 faces from both datasets. Additionally, we manually filtered out low quality meshes. Our final filtered dataset consists of 51k meshes from Objaverse and 5k meshes from ShapeNet. We randomly select 10% of this dataset as the evaluation dataset, with the remaining 90% used as the training set for all our experiments.

Data Processing and Augmentation. Following the strategies of PolyGen [45] and MeshGPT [55], we order faces by their lowest vertex index, then by the next lowest, and so on. Vertices are sorted in ascending order based on their z-y-x coordinates, where z represents the vertical axis. Within each face, we permute the indices to ensure the lowest index comes first. During training, we apply on-the-fly scaling, shifting, and rotation augmentations, normalizing each mesh to a unit bounding box from -0.5 to 0.5 .

5.2 Implementation Details

The encoder and decoder of VQ-VAE both use the encoder of BERT [20], while we choose OPT-350M [74] as our autoregressive transformer architecture. The residual vector quantization [73] depth is set to 3, with a codebook size of 8,192.

Our point encoder is based on the pretrained point encoder from [75], which has been trained on Objaverse and thus can handle general shapes. This point encoder outputs a fixed-length token sequence of 257 tokens, with 256 tokens primarily containing shape information and an additional head token containing semantic information about the shape. We sample 4096 points for each point cloud.

The VQ-VAE is trained on 8 A100 GPUs for 12 hours and the transformer is trained on 8 A100 GPUs for 4 days.

Table 2: Ablation on Noise-Resistant (NR) Decoder for the Quality of Mesh Generation. Please refer to 5.4 for metrics explanation.

Method	CD↓ ($\times 10^{-2}$)	ECD↓ ($\times 10^{-2}$)	NC↑
W/O NR	2.423	6.414	0.883
W/ NR	2.256	6.245	0.902

Table 3: Quantitative evaluation with Marching Cubes [42] and Remesh [4]. Please refer to 5.4 for metrics explanation.

Method	CD↓ ($\times 10^{-2}$)	ECD↓ ($\times 10^{-2}$)	NC↑	#V↓ ($\times 10^3$)	#F↓ ($\times 10^3$)	V_Ratio↓	F_Ratio↓
(a) Marching Cubes	1.532	6.733	0.954	73.22	146.0	440.2	462.2
(b) Remesh (0.005)	2.174	7.813	0.912	127.8	167.9	748.1	534.6
(c) Remesh (0.010)	2.083	7.578	0.929	39.01	41.78	225.4	132.3
(d) Remesh (0.030)	2.915	8.329	0.863	5.848	4.410	34.38	14.05
(e) Remesh (0.050)	4.179	8.138	0.814	2.299	1.538	13.64	4.920
(f) Remesh (0.100)	7.312	10.771	0.748	0.625	0.359	3.735	1.149
(g) <i>MeshAnything</i>	2.256	6.245	0.902	0.172	0.318	0.888	0.871

5.3 Qualitative Experiments

We present more qualitative results of our model combined with other 3D asset production pipelines. As shown in the Fig. 1, Fig. 4 and Fig. 5, MeshAnything effectively generates AMs from various 3D representations. When integrated with different 3D assets production pipelines, our method effectively achieves mesh generation with diverse conditions.

Next, Fig. 5 demonstrates that our model does not simply overfit but understands how to generate meshes with efficient topology that conform to the given shape. To prove this, we use manually-created meshes as ground truth and use their shapes as conditions to test whether our model can generate meshes with comparable topology. To effectively use the ground truth as conditions, we first convert them into dense meshes using Marching Cubes [42] to disrupt their face structure. Then, we sample point clouds with normals from the dense meshes to serve as shape conditions. The experimental results in Fig. 5 show that MeshAnything is capable of generating meshes comparable to or even surpassing those modeled by human artists, exhibiting diverse and strong 3D modeling capabilities.

5.4 Quantitative Experiments

Metrics. We follow the evaluation metric setting of [12]. We quantitatively evaluate mesh quality by uniformly sampling 100K points from the faces of both the ground truth meshes and the predicted meshes, and then computing a set of metrics to assess various aspects of the reconstruction. We report the following metrics: Chamfer Distance (CD) to evaluate the overall quality of a reconstructed mesh; Edge Chamfer Distance (ECD) to assess the preservation of sharp edges by sampling points near sharp edges and corners; and Normal Consistency (NC) to evaluate the quality of the surface normals. Additionally, we report the number of mesh vertices (#V) and the number of mesh faces (#F). We also provide the ratio of the estimated number of vertices to the ground truth number of vertices (#V_Ratio) and the same ratio for faces (#F_Ratio).

Ablations on Noise-Resistant Conditional Decoder. We first perform ablation experiments to verify the effectiveness of the Noise-Resistant Decoder. We begin with a VQ-VAE trained without any noise or conditioning. We then perform ablation between two settings: one where the decoder remains unchanged and unaware of the shape condition, and another where the shape condition is injected into the transformer, as described in Section 4.2. Next, we randomly sample a noise from gumbel distribution and add it to codebook sampling logits during the vector quantization process to simulate the potential low-quality token sequences generated by the transformer. We control the noise level by scaling the added noise.

After training both models for enough epochs, we test their performance to the same level of noise. As shown in Tab. 1, as the intensity of the added noise increases, the Noise-Resistant Decoder with shape condition clearly achieves better reconstruction results. This indicates that the shape condition helps the decoder identify and correct imperfections in the input token sequences.

Next, we verify whether the Noise-Resistant Decoder indeed enhances the transformer’s performance during inference. The test method used dense meshes derived from corrupted GT meshes as the condition for generating new meshes. The generated meshes were then assessed for shape alignment with the conditional shape. As shown in Tab. 2, the model with Noise-Resistant Decoder achieved better results.

Comparison with Marching Cubes and Remesh. Our method is related to various mesh extraction methods [42, 13, 12, 53] since we all convert other 3D representations into meshes. However, it is important to note that previous approaches are reconstruction-like methods that produce dense meshes, while our approach is generative, creating Artist-Created Meshes (AMs) that are significantly more complex to produce than dense meshes. Current metrics can only indicate the quality of shape alignment, which do not effectively reflect the topological advantages of our method. Among the numerous mesh extraction methods, we chose the most representative one, Marching Cubes [42], for comparison.

We compare the results of mesh extraction using Marching Cubes [42] and its combination with Blender remesh [4]. Since our evaluation dataset includes non-watertight meshes, we first extract the signed distance fields (SDF) of all ground truth meshes using [63], which can handle non-watertight meshes. We then apply Marching Cubes with a resolution of 128 on these SDFs. Next, we apply Blender remesh [4] with different voxel sizes to the Marching Cubes results, as both the remesh method and our approach are capable of simplifying topology. Additionally, the Marching Cubes result is used as the shape condition input to MeshAnything to obtain our results.

As shown in Tab. 3, we found that these methods require hundreds of times more faces to achieve results comparable to our method. Furthermore, our approach achieves the best performance in Edge Chamfer Distance (ECD), even with significantly fewer faces than other methods. Comparing (a) and (g), our method lags in Chamfer Distance (CD) and Normal Consistency (NC), mainly due to our method’s inherent failure cases as a generative model, which makes it less robust than Marching Cubes. When comparing with remesh methods, we observe that they incur a high cost to achieve a face count similar to ours. Comparing (f) and (g), we find that even when remesh methods achieve a comparable face count, the number of vertices is still several times higher than ours, indicating that the topology efficiency of remesh methods is far inferior to ours, as they completely ignore the shape characteristics of the 3D assets. Additionally, we surprisingly find that our method can produce results with fewer faces than the ground truth, demonstrating that MeshAnything is not overfitting to the data but instead learns an efficient topology representation, occasionally surpassing the ground truth meshes.

6 Limitations

Our method cannot generate meshes that exceed the maximum face count limit, so it cannot convert large scenes and particularly complex objects into meshes. Additionally, due to its generative nature, our method is not as stable as methods like Marching Cubes [42].

7 Social Impact

Our method points to a promising approach for the automatically generation of Artist-Created Meshes, which has the potential to significantly reduce labor costs in the 3D industry, thereby facilitating advancements in industries such as gaming, film, and the metaverse. However, the reduced cost of obtaining 3D artist-created meshes could also lead to potential criminal activities.

8 Conclusion

In this work, we propose a novel setting for improved mesh extraction and mesh generation, namely Shape-Conditioned Artist-Created Mesh (AM) Generation. Following this setting, we introduce

MeshAnything, a model capable of generating AMs that adhere to given 3D assets. MeshAnything can convert 3D assets in any 3D representation into AMs and thus can be integrated with diverse 3D asset production methods to facilitate their application in the 3D industry. Furthermore, we introduce a noise-resistant decoder architecture to enhance the generation quality, enabling the model to handle low-quality token sequences produced by autoregressive transformers. Lastly, extensive experiments demonstrate the superior performance of our method, highlighting its potential to scale up for 3D industry application and its advantage over previous methods.

Acknowledgments

We thank Yawar Siddiqui for his help with the details about MeshGPT [55], Zibo Zhao for teaching us how to use the Michelangelo [75] encoder, and Ruibo Li for his guidance on point cloud processing.

References

- [1] Antonio Alliegro, Yawar Siddiqui, Tatiana Tommasi, and Matthias Nießner. Polydiff: Generating 3d polygonal meshes with diffusion models. *arXiv preprint arXiv:2312.11417*, 2023.
- [2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.
- [3] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022.
- [4] Blender Development Team. Blender (version 4.1.0) [computer software], 2024. Available from <https://docs.blender.org/manual/en/latest/modeling/modifiers/generate/remesh.html>.
- [5] Jules Bloomenthal. Polygonization of implicit surfaces. *Computer Aided Geometric Design*, 5(4):341–355, 1988.
- [6] Jules Bloomenthal and Chandrajit Bajaj. *Introduction to implicit surfaces*. Morgan Kaufmann, 1997.
- [7] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- [8] Sijin Chen, Xin Chen, Anqi Pang, Xianfang Zeng, Wei Cheng, Yijun Fu, Fukun Yin, Yanru Wang, Zhibin Wang, Chi Zhang, et al. Meshxl: Neural coordinate field for generative 3d foundation models. *arXiv preprint arXiv:2405.20853*, 2024.
- [9] Wenzheng Chen, Huan Ling, Jun Gao, Edward Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. *Advances in neural information processing systems*, 32, 2019.
- [10] Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. *arXiv preprint arXiv:2311.14521*, 2023.
- [11] Yiwen Chen, Chi Zhang, Xiaofeng Yang, Zhongang Cai, Gang Yu, Lei Yang, and Guosheng Lin. It3d: Improved text-to-3d generation with explicit view synthesis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 1237–1244, 2024.
- [12] Zhiqin Chen, Andrea Tagliasacchi, Thomas Funkhouser, and Hao Zhang. Neural dual contouring. *ACM Transactions on Graphics (TOG)*, 41(4):1–13, 2022.
- [13] Zhiqin Chen and Hao Zhang. Neural marching cubes. *ACM Transactions on Graphics (TOG)*, 40(6):1–15, 2021.
- [14] Zilong Chen, Feng Wang, and Huaping Liu. Text-to-3d using gaussian splatting. *arXiv preprint arXiv:2309.16585*, 2023.
- [15] Zilong Chen, Yikai Wang, Feng Wang, Zhengyi Wang, and Huaping Liu. V3d: Video diffusion models are effective 3d generators. *arXiv preprint arXiv:2403.06738*, 2024.

- [16] Evgeni Chernyaev. Marching cubes 33: Construction of topologically correct isosurfaces. Technical report, 1995.
- [17] Morteza Daneshmand, Ahmed Helmi, Egils Avots, Fatemeh Noroozi, Fatih Alisinanoglu, Hasan Sait Arslan, Jelena Gorbova, Rain Eric Haamer, Cagri Ozcinar, and Gholamreza Anbarjafari. 3d scanning: A comprehensive survey. *arXiv preprint arXiv:1801.08863*, 2018.
- [18] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of 10m+ 3d objects. *arXiv preprint arXiv:2307.05663*, 2023.
- [19] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *CVPR*, pages 13142–13153, 2023.
- [20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [21] Akio Doi and Akio Koide. An efficient method of triangulating equi-valued surfaces by using tetrahedral cells. *IEICE TRANSACTIONS on Information and Systems*, 74(1):214–224, 1991.
- [22] Jiemin Fang, Junjie Wang, Xiaopeng Zhang, Lingxi Xie, and Qi Tian. Gaussianeditor: Editing 3d gaussians delicately with text instructions. *arXiv preprint arXiv:2311.16037*, 2023.
- [23] Jun Gao, Wenzheng Chen, Tommy Xiang, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Learning deformable tetrahedral meshes for 3d reconstruction. *Advances in neural information processing systems*, 33:9936–9947, 2020.
- [24] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from images. *NeurIPS*, 35:31841–31854, 2022.
- [25] Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen, Peng Gao, Xianzhi Li, Hongsheng Li, et al. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint arXiv:2309.00615*, 2023.
- [26] Abid Haleem and Mohd Javaid. 3d scanning applications in medical field: a literature-based review. *Clinical Epidemiology and Global Health*, 7(2):199–210, 2019.
- [27] Abid Haleem, Mohd Javaid, Ravi Pratap Singh, Shanay Rab, Rajiv Suman, Lalit Kumar, and Ibrahim Haleem Khan. Exploring the potential of 3d scanning in industry 4.0: An overview. *International Journal of Cognitive Computing in Engineering*, 3:161–171, 2022.
- [28] Rana Hanocka, Gal Metzer, Raja Giryes, and Daniel Cohen-Or. Point2mesh: A self-prior for deformable meshes. *arXiv preprint arXiv:2005.11084*, 2020.
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [30] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint arXiv:2311.04400*, 2023.
- [31] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for geometrically accurate radiance fields. *arXiv preprint arXiv:2403.17888*, 2024.
- [32] Tao Ju, Frank Losasso, Scott Schaefer, and Joe Warren. Dual contouring of hermite data. In *Proceedings of the 29th annual conference on Computer graphics and interactive techniques*, pages 339–346, 2002.
- [33] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3907–3916, 2018.
- [34] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023.

- [35] Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model. *arXiv preprint arXiv:2311.06214*, 2023.
- [36] Yiyi Liao, Simon Donne, and Andreas Geiger. Deep marching cubes: Learning explicit surface representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2916–2925, 2018.
- [37] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [38] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *Advances in Neural Information Processing Systems*, 36, 2024.
- [39] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. <https://arxiv.org/abs/2303.11328>, 2023.
- [40] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. *arXiv preprint arXiv:2303.11328*, 2023.
- [41] Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma, Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. Wonder3d: Single image to 3d using cross-domain diffusion. *arXiv preprint arXiv:2310.15008*, 2023.
- [42] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987.
- [43] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *Seminal graphics: pioneering efforts that shaped the field*, pages 347–353. 1998.
- [44] 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. In *ECCV*, 2020.
- [45] Charlie Nash, Yaroslav Ganin, SM Ali Eslami, and Peter Battaglia. Polygen: An autoregressive generative model of 3d meshes. In *International conference on machine learning*, pages 7220–7229. PMLR, 2020.
- [46] Gregory M Nielson. Dual marching cubes. In *IEEE visualization 2004*, pages 489–496. IEEE, 2004.
- [47] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- [48] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR 2017*, pages 652–660, 2017.
- [49] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017.
- [50] Scott Schaefer, Tao Ju, and Joe Warren. Manifold dual contouring. *IEEE Transactions on Visualization and Computer Graphics*, 13(3):610–619, 2007.
- [51] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. *Advances in Neural Information Processing Systems*, 34:6087–6101, 2021.
- [52] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [53] Tianchang Shen, Jacob Munkberg, Jon Hasselgren, Kangxue Yin, Zian Wang, Wenzheng Chen, Zan Gojcic, Sanja Fidler, Nicholas Sharp, and Jun Gao. Flexible isosurface extraction for gradient-based mesh optimization. *ACM Transactions on Graphics (TOG)*, 42(4):1–16, 2023.

- [54] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023.
- [55] Yawar Siddiqui, Antonio Alliegro, Alexey Artemov, Tatiana Tommasi, Daniele Sirigatti, Vladislav Rosov, Angela Dai, and Matthias Nießner. Meshgpt: Generating triangle meshes with decoder-only transformers. *arXiv preprint arXiv:2311.15475*, 2023.
- [56] Jingxiang Sun, Bo Zhang, Ruizhi Shao, Lizhen Wang, Wen Liu, Zhenda Xie, and Yebin Liu. Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior. *arXiv preprint arXiv:2310.16818*, 2023.
- [57] Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm: Large multi-view gaussian model for high-resolution 3d content creation. *arXiv preprint arXiv:2402.05054*, 2024.
- [58] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.
- [59] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. *arXiv preprint arXiv:2303.14184*, 2023.
- [60] Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, Adam Letts, Yangguang Li, Ding Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. Triposr: Fast 3d object reconstruction from a single image. *arXiv preprint arXiv:2403.02151*, 2024.
- [61] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- [62] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [63] Peng-Shuai Wang, Yang Liu, and Xin Tong. Dual octree graph networks for learning adaptive volumetric shape representations. *ACM Transactions on Graphics (TOG)*, 41(4):1–15, 2022.
- [64] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolicdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv preprint arXiv:2305.16213*, 2023.
- [65] Zhengyi Wang, Yikai Wang, Yifei Chen, Chendong Xiang, Shuo Chen, Dajiang Yu, Chongxuan Li, Hang Su, and Jun Zhu. Crm: Single image to 3d textured mesh with convolutional reconstruction model. *arXiv preprint arXiv:2403.05034*, 2024.
- [66] Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan, Valentin Deschaintre, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Meshlrn: Large reconstruction model for high-quality mesh. *arXiv preprint arXiv:2404.12385*, 2024.
- [67] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pages 6861–6871. PMLR, 2019.
- [68] Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and S Yu Philip. Multimodal large language models: A survey. In *2023 IEEE International Conference on Big Data (BigData)*, pages 2247–2256. IEEE, 2023.
- [69] Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. Instantmesh: Efficient 3d mesh generation from a single image with sparse-view large reconstruction models. *arXiv preprint arXiv:2404.07191*, 2024.
- [70] Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. Pointllm: Empowering large language models to understand point clouds. *arXiv preprint arXiv:2308.16911*, 2023.
- [71] Xiaofeng Yang, Yiwen Chen, Cheng Chen, Chi Zhang, Yi Xu, Xulei Yang, Fayao Liu, and Guosheng Lin. Learn to optimize denoising scores for 3d generation: A unified and improved diffusion prior on nerf and 3d gaussian splatting. *arXiv preprint arXiv:2312.04820*, 2023.
- [72] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *CVPR*, pages 4578–4587, 2021.

- [73] Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. Soundstream: An end-to-end neural audio codec. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:495–507, 2021.
- [74] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [75] Zibo Zhao, Wen Liu, Xin Chen, Xianfang Zeng, Rui Wang, Pei Cheng, Bin Fu, Tao Chen, Gang Yu, and Shenghua Gao. Michelangelo: Conditional 3d shape generation based on shape-image-text aligned latent representation. *Advances in Neural Information Processing Systems*, 36, 2024.