

MESHANYTHING: ARTIST-CREATED MESH GENERATION WITH AUTOREGRESSIVE TRANSFORMERS

Anonymous authors

Paper under double-blind review

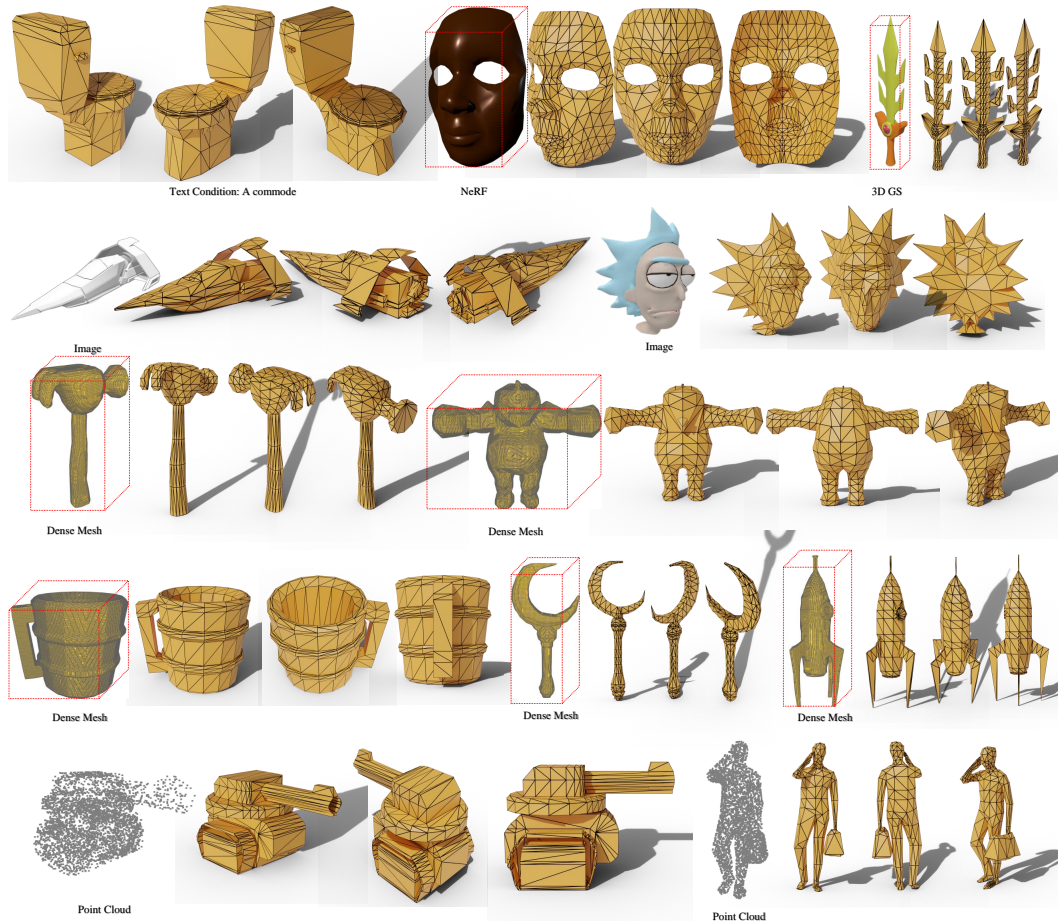


Figure 1: **MeshAnything** converts different 3D representation into **Artist-Created Meshes (AMs)**, i.e., meshes created by human artists. It can be combined with various 3D asset production pipelines, such as 3D reconstruction and generation, to transform their results into AMs that can be seamlessly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

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.

054 By converting 3D assets in any 3D representation into AMs, MeshAnything can
055 be integrated with various 3D asset production methods, thereby enhancing their
056 application across the 3D industry. The architecture of MeshAnything comprises
057 a VQ-VAE and a shape-conditioned decoder-only transformer. We first learn a
058 mesh vocabulary using the VQ-VAE, then train the shape-conditioned decoder-
059 only transformer on this vocabulary for shape-conditioned autoregressive mesh
060 generation. Our extensive experiments show that our method generates AMs with
061 hundreds of times fewer faces, significantly improving storage, rendering, and
062 simulation efficiencies, while achieving precision comparable to previous meth-
063 ods.
064

065 1 INTRODUCTION

066

067 In recent years, the 3D community has experienced rapid advancements, with a variety of methods
068 developed for automatically producing high-quality 3D assets. These methods, including 3D recon-
069 struction (Mildenhall et al., 2020; Yu et al., 2021; Barron et al., 2021; 2022; Kerbl et al., 2023b;
070 Huang et al., 2024), 3D generation (Poole et al., 2023; Liu et al., 2023a; Wang et al., 2023; Long
071 et al., 2023; Sun et al., 2023; Hong et al., 2023; Tang et al., 2024; Xu et al., 2024; Wei et al., 2024),
072 and scanning (Daneshmand et al., 2018; Haleem & Javaid, 2019; Haleem et al., 2022), can produce
073 3D assets with shape and color quality comparable to manually created ones. The success of these
074 methods reveals the potential to replace manually created 3D models with automatically produced
075 ones in the 3D industry, including applications in games, movies, and the metaverse, significantly
076 reducing time and labor costs.

077 However, this potential remains largely unrealized because the current 3D industry predominantly
078 relies on mesh-based pipelines for their superior efficiency and controllability, while methods
079 for producing 3D assets typically use alternative 3D representations to achieve optimal results
080 across various scenarios. Therefore, substantial efforts (Lorensen & Cline, 1987; Chernyaev, 1995;
081 Lorensen & Cline, 1998; Shen et al., 2021b; Chen et al., 2022; Shen et al., 2023) are devoted to
082 converting other 3D representations into meshes and have achieved some success. Meshes produced
083 by these methods approximate the shape quality of those created by human artists, which we refer
084 to as Artist-Created Meshes (AMs), but they still fall short in addressing the aforementioned issues.

085 This is because all meshes produced by these methods (Lorensen & Cline, 1987; Chernyaev, 1995;
086 Lorensen & Cline, 1998; Shen et al., 2021b; Chen et al., 2022; Shen et al., 2023) exhibit significantly
087 poorer topology quality compared to AMs. As shown in Fig. 2, these methods rely on dense faces to
088 reconstruct 3D shapes, completely ignoring geometric characteristics. Using these meshes in the 3D
089 industry leads to three significant problems: First, converted meshes typically contain several orders
090 of magnitude more faces compared to AMs, leading to significant inefficiencies in storage, render-
091 ing, and simulation. Moreover, the converted meshes complicate post-processing and downstream
092 tasks in the 3D pipeline. They significantly increase the challenge for human artists in optimizing
093 these meshes due to their chaotic and inefficient topologies. Finally, previous methods struggle to
094 represent sharp edges and flat surfaces, resulting in oversmoothing and bumpy artifacts as shown in
095 Fig. 2.

095 In this work, we aim to solve the aforementioned issues to facilitate the application of automatically
096 generated 3D assets in the 3D industry. As mentioned earlier, all previous methods (Lorensen &
097 Cline, 1987; Chernyaev, 1995; Lorensen & Cline, 1998; Shen et al., 2021b; Chen et al., 2022; Shen
098 et al., 2023) extract 3D meshes with excessively dense faces in a reconstruction manner, which in-
099 herently cannot solve these issues. Therefore, we diverge from previous approaches by formulating
100 mesh extraction as a generation problem for the first time: we teach models to generate Artist-
101 Created Meshes (AMs) that are aligned with the given 3D assets. The meshes generated by our
102 methods mimic the shape and topology quality of those created by human artists. Consequently, our
103 setting, namely Shape-Conditioned AM Generation, is fundamentally free from all previous issues,
104 enabling seamless integration of the generated results into the 3D industry pipeline.

105 However, training such a model presents significant challenges. The first challenge is constructing
106 the dataset, as we need paired shape conditions and Artist-Created Meshes (AMs) for model train-
107 ing. 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

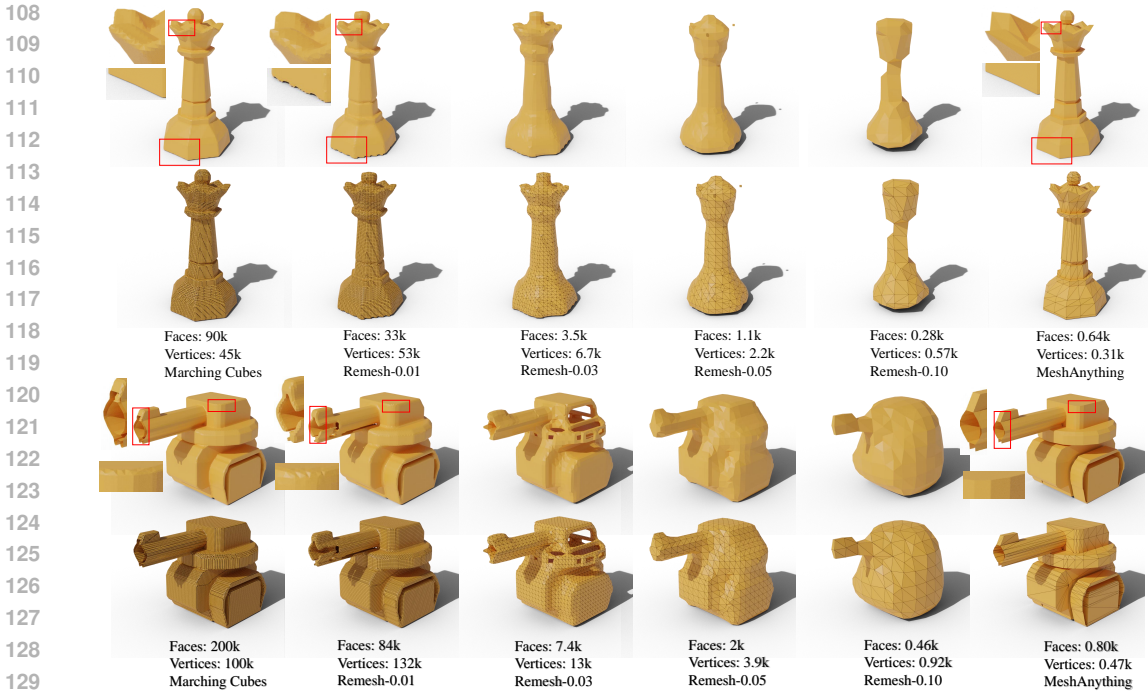


Figure 2: **Comparison with Marching Cubes Lorensen & Cline (1987) and Remesh Blender Development Team (2024).** We apply Marching Cubes and MeshAnything to ground truth shapes and then apply remeshing to the Marching Cubes results with different voxel sizes. Existing methods extract meshes in a reconstruction 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.

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 (Qi et al., 2017a;b; Zhao et al., 2024).

We filter out high-quality AMs from Objaverse (Deitke et al., 2023b;a) and ShapeNet (Chang et al., 2015). 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 (Wang et al., 2022), convert it into a coarser mesh using (Lorensen & Cline, 1987), and then sample point clouds from this coarse mesh to narrow the domain gap in shape conditions between inference and training.

Following (Siddiqui et al., 2023), we use a VQ-VAE (Van Den Oord et al., 2017) to learn a mesh vocabulary and train a decoder-only transformer (Vaswani et al., 2017) on this vocabulary for mesh generation. To inject shape condition, we draw inspiration from the recent success of multimodal large language models (MLLM) (Wu et al., 2023; Liu et al., 2024b), 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 (Zhao et al., 2024) 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 (Siddiqui et al., 2023).

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 (Van Den Oord et al.,

2017) 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 as follows:

- 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.
- Extensive experiments demonstrate that Shape-Conditioned Mesh Generation is a more suitable setting for mesh generation, and MeshAnything significantly surpasses previous mesh generation methods.

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 (Shen et al., 2023), we categorize these methods into two main types: Isosurface Extraction (Lorensen & Cline, 1987; Bloomenthal, 1988; Chernyaev, 1995; Bloomenthal & Bajaj, 1997; Lorensen & Cline, 1998; Chen et al., 2022) and Gradient-Based Mesh Optimization (Chen et al., 2019; Gao et al., 2020; Hanocka et al., 2020; Kato et al., 2018; Shen et al., 2021a; Liao et al., 2018; Shen et al., 2023).

Traditional isosurface extraction methods (Lorensen & Cline, 1987; 1998; Chernyaev, 1995; Doi & Koide, 1991; Ju et al., 2002; Schaefer et al., 2007; Chen & Zhang, 2021; Chen et al., 2022) 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 (Lorensen & Cline, 1987). 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. Recently, (Chen & Zhang, 2021) and (Chen et al., 2022) introduce data-driven methods 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 (Chen et al., 2019; Gao et al., 2020; Hanocka et al., 2020; Kato et al., 2018; Shen et al., 2021a; Liao et al., 2018; Shen et al., 2023). 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. (Shen et al., 2021a) employs a differentiable Marching Tetrahedra layer for mesh extraction. Similar to (Shen et al., 2021a), (Shen et al., 2023) iteratively optimizes a 3D surface mesh by representing it as the isosurface of a scalar field.

216 However, these approaches fundamentally differ from ours. They ignore the characteristics of the
 217 shape and inherently cannot produce meshes with efficient topology. In contrast, MeshAnything
 218 formulates mesh extraction as a generation problem for the first time, aiming to mimic human artists
 219 in mesh extraction and thereby generating Artist-Created Meshes (AMs) with hundreds of times
 220 fewer faces.

222 2.2 3D MESH GENERATIONS

224 3D mesh generation can be mainly divided into two categories: generating dense meshes similar to
 225 those produced by previous mesh extraction methods, and generating Artist-Created Meshes (AMs).

226 The former category is currently the mainstream research focus. Methods such as (Gao et al., 2022;
 227 Wei et al., 2024; Xu et al., 2024) directly generate meshes in a feed-forward manner, but because
 228 they produce dense meshes with low-quality topology similar to previous mesh extraction methods,
 229 they still encounter the same issues when applied in the 3D industry.

230 Notably, numerous 3D generation methods (Poole et al., 2023; Tang et al., 2023b; Wang et al., 2023;
 231 Chen et al., 2024b; Tang et al., 2023a; Yang et al., 2023; Hong et al., 2023; Fang et al., 2023; Chen
 232 et al., 2023a; Liu et al., 2024c; Shi et al., 2023; Li et al., 2023; Chen et al., 2023b; 2024c; Tang et al.,
 233 2024; Wang et al., 2024; Tochilkin et al., 2024; Liu et al., 2024a; 2023c; Zhang et al., 2024) can also
 234 produce meshes. These methods first generate 3D assets and then convert them to dense meshes
 235 using mesh extraction methods like (Lorensen & Cline, 1987). Consequently, they face challenges
 236 when applied to the 3D industry due to their inefficient topology.

237 Recently, several works have focused on the second category: generating Artist-Created
 238 Meshes(AMs) (Nash et al., 2020; Alliegro et al., 2023; Siddiqui et al., 2023; Chen et al., 2024a).
 239 Although our approach also focuses on AM generation, it fundamentally differs from these meth-
 240 ods. Since they lack shape conditioning, these methods must simultaneously learn the complex 3D
 241 shape distribution—which typically alone requires extensive training (Hong et al., 2023; Tang et al.,
 242 2024)—and the topology distribution of AMs, leading to very challenging training processes. In
 243 contrast, our methods eliminate the challenge of learning the shape distribution, allowing the model
 244 to focus on learning the topology distribution. This not only significantly reduces training costs but
 245 also enhances the model’s application value.

246 Among these methods, the most relevant to ours is MeshGPT (Siddiqui et al., 2023), as we follow
 247 its architecture. (Siddiqui et al., 2023) introduced a combination of a VQ-VAE (Van Den Oord et al.,
 248 2017) and an autoregressive transformer architecture. It first learns a mesh vocabulary with the
 249 VQ-VAE and then trains the transformer on the learned vocabulary for mesh generation. However,
 250 MeshGPT’s results are limited to several categories in ShapeNet. MeshGPT requires a training GPU
 251 hours similar to ours, but our method can generalize to unlimited categories in Objaverse. As shown
 252 in Fig. 3, this is largely due to the difference in target complexity caused by MeshGPT needing to
 253 additionally learn the complex 3D shape distribution.

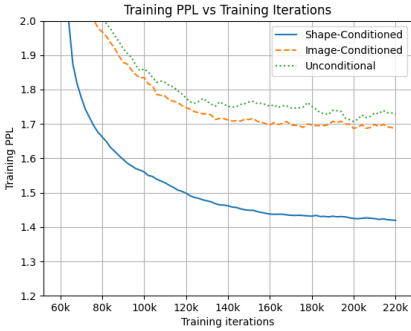
255 3 SHAPE-CONDITIONED AM GENERATION

257 In this section, we first introduce the formal formulation for Shape-Conditioned AM Generation and
 258 compare it with previous mesh generation settings (Nash et al., 2020; Siddiqui et al., 2023; Alliegro
 259 et al., 2023). We show that it can achieve better performance and a broader range of applications
 260 compared to the settings in previous mesh generation methods, with significantly less training effort.

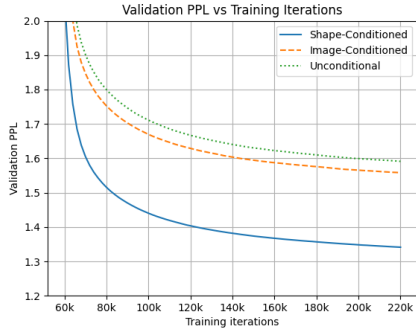
262 Shape-Conditioned AM Generation targets to estimate a conditional distribution $p(\mathcal{M}|\mathcal{S})$. In this
 263 formula, \mathcal{M} refers to the Artist-Created Mesh (AM), i.e., the mesh manually modeled by human
 264 artists. \mathcal{S} refers to the 3D shape information that indicates the 3D shape to which \mathcal{M} should align.
 265 The input form of \mathcal{S} can be diverse, such as voxels or point clouds. Therefore, this versatility allows
 266 our method to be integrated with any 3D pipeline that outputs \mathcal{S} , such as 3D reconstruction (Milden-
 267 hall et al., 2020; Kerbl et al., 2023b), generation (Poole et al., 2023; Hong et al., 2023), and scanning,
 268 making these methods more efficient for the 3D industry.

269 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,

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323



(a) Training Perplexity (PPL)



(b) Validation Perplexity (PPL)

Figure 3: **Training and validation perplexity (PPL) for the mesh generation model under different input conditions.** All models are trained with the same settings as detailed in Section 5.2. The training and validation PPL of shape-conditioned mesh generation is significantly lower than that of unconditional and image-conditioned mesh generation. This indicates that the training burden of shape-conditioned mesh generation is much lower since it avoids learning the complex 3D shape distribution.

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:

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

According to the chain rule, we have:

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

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:

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

Combining 1, 2 and 3:

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

in which $p(\mathcal{M}|\mathcal{S})$ is the focus of our shape-conditioned mesh generation. As shown in Fig. 3, estimating $p(\mathcal{M}|\mathcal{S})$ is much more simpler than $p(\mathcal{M}|\mathcal{C})$, proving that our setting is much easier to train than settings in prvious methods.

As for $p(\mathcal{S}|\mathcal{C})$, In the 3D community, numerous large models (Team, 2024; Tang et al., 2024; Xu et al., 2024; Siddiqui et al., 2023) aim to estimate using various 3D representations and demonstrate excellent results. Besides, some single scene 3D asset production methods (Mildenhall et al., 2020; Kerbl et al., 2023b; Barron et al., 2021; 2022; Poole et al., 2023; Liu et al., 2023b; Sun et al., 2023) 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:

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

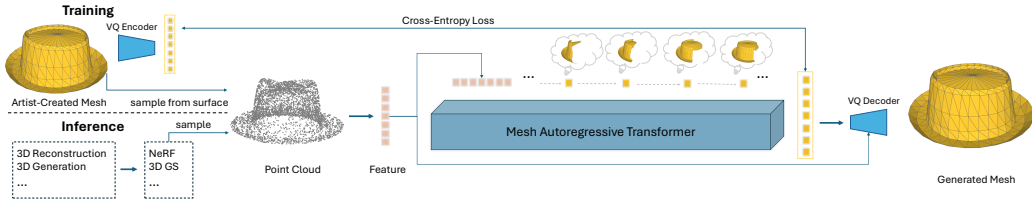


Figure 4: **Pipeline Overview.** We introduce MeshAnything, an autoregressive transformer capable of generating Artist-Created Meshes that adhere to given 3D shapes. During training, we inject point clouds features into a decoder-only transformer and supervise it using token sequences derived from the Artist-Created meshes. After training, MeshAnything takes point clouds sampled from various 3D representations as input and generates aligned Artist-Created meshes.

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 (Mildenhall et al., 2020; Kerbl et al., 2023b; Hong et al., 2023; Poole et al., 2023; Tang et al., 2024).
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 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 (Wang et al., 2022), then convert it into a relatively coarse mesh using Marching Cubes (Lorensen & Cline, 1987) 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 (Lorensen & Cline, 1987) 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 (Siddiqui et al., 2023), we first train a VQ-VAE (Van Den Oord et al., 2017) to learn a vocabulary of geometric embeddings for better transformer (Vaswani et al., 2017) learning. Different to MeshGPT, which uses graph convolutional networks (Wu et al., 2019) and ResNet (He et al., 2016) 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

378 as a sequence of triangle faces:

$$379 \quad \mathcal{M} := (f_1, f_2, f_3, \dots, f_N), \quad (5)$$

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

$$384 \quad \mathcal{Z} = (z_1, z_2, \dots, z_N) = E(\mathcal{M}), \quad (6)$$

385 where z_i is the feature vector for f_i .

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

$$388 \quad \mathcal{T} = RQ(\mathcal{Z}; \mathcal{B}) \quad (7)$$

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

$$393 \quad \hat{\mathcal{M}} = D(\mathcal{Z}) \quad (8)$$

394 The VQ-VAE is trained end-to-end with cross-entropy loss on the predicted vertex coordinate logits
395 and the commitment loss of vector quantization (Van Den Oord et al., 2017). After the training of
396 VQ-VAE, the encoder-decoder of VQ-VAE is treated as a tokenizer and detokenizer for autoregres-
397 sive transformer training.

399 However, as shown in Fig. 7, there are possible imperfections in the generation results. To address
400 this issue, given our setting of Shape-Conditioned AM Generation, the VQ-VAE decoder can also
401 take the shape condition as input. Small imperfections in the token sequences generated by the
402 transformer can potentially be corrected by a shape-aware decoder. Therefore, after completing the
403 vanilla VQ-VAE training, we add an additional decoder fine-tuning stage, where we inject the shape
404 information into the transformer decoder. Then we add random Gumbel noise to the codebook
405 sampling logits to simulate the potential imperfections in the token sequences generated by the
406 transformer during inference. The decoder is then updated independently with the same cross-
407 entropy loss to train it to produce refined meshes even when facing imperfect token sequences. Our
408 experiments in Tab. 3 and Tab. 4 show that our method effectively enhances the decoder’s noise
409 resistance and mesh generation quality.

410 4.3 SHAPE-CONDITIONED AUTOREGRESSIVE TRANSFORMER

411 To add shape condition to the transformer, inspired by the success of multimodal large language
412 models (Wu et al., 2023; Liu et al., 2024b; Xu et al., 2023; Guo et al., 2023), we first encode the
413 point cloud into a fixed-length token sequence with a point cloud encoder \mathcal{P} and then concatenate it
414 to the front of the embedding sequence from \mathcal{T} VQ-VAE as the final input embedding sequence for
415 the transformer:

$$416 \quad \mathcal{T}' = \text{concat}(\mathcal{P}(S), \mathcal{T}) \quad (9)$$

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

418 We borrow a pretrained point encoder from (Zhao et al., 2024) and add a linear projection layer
419 to project its output feature to the same latent space as \mathcal{T} . During training, the original point en-
420 coder from (Zhao et al., 2024) is frozen; we only update the newly added projection layer and the
421 autoregressive transformer with cross-entropy loss.

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

$$424 \quad \hat{\mathcal{M}} = D(\hat{\mathcal{T}}) \quad (10)$$

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

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

Table 1: Comparison of Mesh Generation Methods. As shown in the left table, compared to the baseline Artist-Created Mesh Generation method, the meshes generated by MeshAnything are better aligned with human preferences. In the right table, we compare MeshAnything with mesh extraction baselines, and it received the most votes. For detailed settings, please refer to Section 5.4.

Method	Shape \uparrow	Topology \uparrow	Method	Shape \uparrow	Topology \uparrow
PolyGen	12.7%	11.1%	MarchingCubes	38.1%	10.2%
MeshGPT	24.1%	28.2%	Shape As Points	17.3%	6.2%
MeshAnything	63.2%	60.7%	MeshAnything	44.6%	83.6%

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 (Deitke et al., 2023b) and ShapeNet (Chang et al., 2015), 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 (Nash et al., 2020) and MeshGPT (Siddiqui et al., 2023), 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 (Devlin et al., 2018), while we choose OPT-350M (Zhang et al., 2022) as our autoregressive transformer architecture. The residual vector quantization (Zeghidour et al., 2021) depth is set to 3, with a codebook size of 8,192.

Our point encoder is based on the pretrained point encoder from (Zhao et al., 2024), 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 training batch size for both the VQ-VAE and the transformer is set to 8 per GPU. The VQ-VAE is trained on 8 A100 GPUs for 12 hours, after which we separately finetune the decoder part of the VQ-VAE into a noise-resistant decoder, as detailed in Section 4.2. Following this, the transformer is trained on 8 A100 GPUs for 4 days.

Table 2: **Quantitative Comparisons with Prior Arts on Objaverse.** MeshAnything significantly outperforms prior methods across all metrics. MMD, KID are scaled by 10^3 .

Method	COV \uparrow	MMD \downarrow	1-NNA \downarrow	FID \downarrow	KID \downarrow
PolyGen	23.2	6.22	88.2	48.8	27.7
MeshGPT	41.7	3.83	67.3	25.1	6.11
MeshAnything	53.1	2.72	55.7	14.5	1.89

5.3 QUALITATIVE EXPERIMENTS

As shown in Fig. 1, MeshAnything effectively generates AMs from various 3D representations. In our experiments, we use Rodin (Team, 2024) as the text-to-3D and image-to-3D method, and employ (Mildenhall et al., 2020) and (Kerbl et al., 2023a) as the 3D reconstruction pipeline to obtain the corresponding NeRF and Gaussian Splatting models. For additional qualitative results, please refer to A.2 combined with other 3D asset production pipelines.

5.4 QUANTITATIVE EXPERIMENTS

From the generative model perspective, MeshAnything is a shape-conditioned mesh generation model. From the mesh extraction perspective, it extracts artist-created meshes from point clouds. Consequently, we compare MeshAnything with both types of methods. Additional experiments can be found in Appendix Section A.2.

User Study. As shown in Tab. 1, we conducted two user studies, comparing with mesh generation baselines (Nash et al., 2020; Siddiqui et al., 2023) and mesh extraction baselines (Lorensen & Cline, 1987; Peng et al., 2021), respectively. The majority of participants in our user study were researchers from the 3D field, with a smaller portion being 3D industry practitioners. The mesh generation baselines are trained on ShapeNet, and to ensure a fair comparison, we retrained them on Objaverse using the same transformer model as MeshAnything. Since the mesh generation baselines are all unconditional mesh generation methods, whereas MeshAnything is a shape-conditioned mesh generation method, we sampled shapes randomly from the evaluation set of Objaverse as inputs for MeshAnything, while for the baseline methods, we performed random sampling directly.

In the mesh extraction baseline, since our method can also be viewed as a point cloud to mesh approach, we included (Peng et al., 2021), a point cloud to mesh method, as a baseline. Additionally, we optimized the results from the mesh extraction baseline using the Blender remesh method (Blender Development Team, 2024) to simplify the topology.

We collected 30 results from each method and asked users to vote for the best one in terms of shape quality and topology quality. A total of 41 users participated, providing 1,230 valid comparisons. Both user studies demonstrated the superiority of our method. The only difference between the retrained MeshGPT and MeshAnything is whether they are shape-conditioned, further proving the advantages of the shape-conditioned mesh generation setting.

Metrics. We follow the metric setting of (Chen et al., 2022; Siddiqui et al., 2023). We detail this setting in Appendix Section. A.1.

Comparison with Mesh Generation Pipelines. We use the same retrained models from the user study for comparison. As shown in Tab. 2, MeshAnything significantly outperforms prior methods (Nash et al., 2020; Siddiqui et al., 2023), indicating that it’s superior in both the shape and topology quality. Since the only difference between the retrained MeshGPT and MeshAnything is the inclusion of shape conditioning, the superior performance of MeshAnything further demonstrates that Shape-Conditioned Mesh Generation is a more suitable setting for mesh generation.

6 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.

REFERENCES

- 540
541
542 Antonio Alliegro, Yawar Siddiqui, Tatiana Tommasi, and Matthias Nießner. Polydiff: Generating
543 3d polygonal meshes with diffusion models. *arXiv preprint arXiv:2312.11417*, 2023.
- 544
545 Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and
546 Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields.
547 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5855–5864,
548 2021.
- 549
550 Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf
360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022.
- 551
552 Blender Development Team. Blender (version 4.1.0) [computer software], 2024. Available
553 from [https://docs.blender.org/manual/en/latest/modeling/modifiers/
554 generate/remesh.html](https://docs.blender.org/manual/en/latest/modeling/modifiers/generate/remesh.html).
- 555
556 Jules Bloomenthal. Polygonization of implicit surfaces. *Computer Aided Geometric Design*, 5(4):
341–355, 1988.
- 557
558 Jules Bloomenthal and Chandrajit Bajaj. *Introduction to implicit surfaces*. Morgan Kaufmann, 1997.
- 559
560 Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li,
561 Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d
model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- 562
563 Sijin Chen, Xin Chen, Anqi Pang, Xianfang Zeng, Wei Cheng, Yijun Fu, Fukun Yin, Yanru Wang,
564 Zhibin Wang, Chi Zhang, et al. Meshxl: Neural coordinate field for generative 3d foundation
565 models. *arXiv preprint arXiv:2405.20853*, 2024a.
- 566
567 Wenzheng Chen, Huan Ling, Jun Gao, Edward Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja
568 Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. *Ad-
vances in neural information processing systems*, 32, 2019.
- 569
570 Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei
571 Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with
572 gaussian splatting. *arXiv preprint arXiv:2311.14521*, 2023a.
- 573
574 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
575 Conference on Artificial Intelligence*, volume 38, pp. 1237–1244, 2024b.
- 576
577 Zhiqin Chen and Hao Zhang. Neural marching cubes. *ACM Transactions on Graphics (TOG)*, 40
578 (6):1–15, 2021.
- 579
580 Zhiqin Chen, Andrea Tagliasacchi, Thomas Funkhouser, and Hao Zhang. Neural dual contouring.
ACM Transactions on Graphics (TOG), 41(4):1–13, 2022.
- 581
582 Zilong Chen, Feng Wang, and Huaping Liu. Text-to-3d using gaussian splatting. *arXiv preprint
583 arXiv:2309.16585*, 2023b.
- 584
585 Zilong Chen, Yikai Wang, Feng Wang, Zhengyi Wang, and Huaping Liu. V3d: Video diffusion
586 models are effective 3d generators. *arXiv preprint arXiv:2403.06738*, 2024c.
- 587
588 Evgeni Chernyaev. Marching cubes 33: Construction of topologically correct isosurfaces. Technical
report, 1995.
- 589
590 Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation,
591 Stichting Blender Foundation, Amsterdam, 2018. URL <http://www.blender.org>.
- 592
593 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.

- 594 Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan
595 Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of
596 10m+ 3d objects. *arXiv preprint arXiv:2307.05663*, 2023a.
- 597 Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig
598 Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of anno-
599 tated 3d objects. In *CVPR*, pp. 13142–13153, 2023b.
- 600 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
601 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 602
603 Akio Doi and Akio Koide. An efficient method of triangulating equi-valued surfaces by using
604 tetrahedral cells. *IEICE TRANSACTIONS on Information and Systems*, 74(1):214–224, 1991.
- 605
606 Jiemin Fang, Junjie Wang, Xiaopeng Zhang, Lingxi Xie, and Qi Tian. Gaussianeditor: Editing 3d
607 gaussians delicately with text instructions. *arXiv preprint arXiv:2311.16037*, 2023.
- 608
609 Jun Gao, Wenzheng Chen, Tommy Xiang, Alec Jacobson, Morgan McGuire, and Sanja Fidler.
610 Learning deformable tetrahedral meshes for 3d reconstruction. *Advances in neural information
611 processing systems*, 33:9936–9947, 2020.
- 612
613 Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan
614 Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned
615 from images. *NeurIPS*, 35:31841–31854, 2022.
- 616
617 Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen,
618 Peng Gao, Xianzhi Li, Hongsheng Li, et al. Point-bind & point-llm: Aligning point cloud
619 with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint
620 arXiv:2309.00615*, 2023.
- 621
622 Abid Haleem and Mohd Javaid. 3d scanning applications in medical field: a literature-based review.
623 *Clinical Epidemiology and Global Health*, 7(2):199–210, 2019.
- 624
625 Abid Haleem, Mohd Javaid, Ravi Pratap Singh, Shanay Rab, Rajiv Suman, Lalit Kumar, and
626 Ibrahim Haleem Khan. Exploring the potential of 3d scanning in industry 4.0: An overview.
627 *International Journal of Cognitive Computing in Engineering*, 3:161–171, 2022.
- 628
629 Rana Hanocka, Gal Metzer, Raja Giryes, and Daniel Cohen-Or. Point2mesh: A self-prior for de-
630 formable meshes. *arXiv preprint arXiv:2005.11084*, 2020.
- 631
632 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
633 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
634 770–778, 2016.
- 635
636 Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli,
637 Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint
638 arXiv:2311.04400*, 2023.
- 639
640 Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting
641 for geometrically accurate radiance fields. *arXiv preprint arXiv:2403.17888*, 2024.
- 642
643 Tao Ju, Frank Losasso, Scott Schaefer, and Joe Warren. Dual contouring of hermite data. In *Pro-
644 ceedings of the 29th annual conference on Computer graphics and interactive techniques*, pp.
645 339–346, 2002.
- 646
647 Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *Proceedings of
648 the IEEE conference on computer vision and pattern recognition*, pp. 3907–3916, 2018.
- 649
650 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splat-
651 ting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14,
652 2023a.
- 653
654 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splat-
655 ting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023b.
656 URL <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>.

- 648 Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan
649 Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3d: Fast text-to-3d with sparse-view gen-
650 eration and large reconstruction model. *arXiv preprint arXiv:2311.06214*, 2023.
- 651 Yiyi Liao, Simon Donne, and Andreas Geiger. Deep marching cubes: Learning explicit surface
652 representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recog-
653 nition*, pp. 2916–2925, 2018.
- 654 Anran Liu, Cheng Lin, Yuan Liu, Xiaoxiao Long, Zhiyang Dou, Hao-Xiang Guo, Ping Luo, and
655 Wenping Wang. Part123: Part-aware 3d reconstruction from a single-view image, 2024a. URL
656 <https://arxiv.org/abs/2405.16888>.
- 657 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances
658 in neural information processing systems*, 36, 2024b.
- 659 Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. One-
660 2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *Advances in
661 Neural Information Processing Systems*, 36, 2024c.
- 662 Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick.
663 Zero-1-to-3: Zero-shot one image to 3d object. <https://arxiv.org/abs/2303.11328>, 2023a.
- 664 Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick.
665 Zero-1-to-3: Zero-shot one image to 3d object. *arXiv preprint arXiv:2303.11328*, 2023b.
- 666 Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang.
667 Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint
668 arXiv:2309.03453*, 2023c.
- 669 Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma,
670 Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. Wonder3d: Single image to 3d
671 using cross-domain diffusion. *arXiv preprint arXiv:2310.15008*, 2023.
- 672 William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction
673 algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987.
- 674 William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction
675 algorithm. In *Seminal graphics: pioneering efforts that shaped the field*, pp. 347–353. 1998.
- 676 Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and
677 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- 678 Charlie Nash, Yaroslav Ganin, SM Ali Eslami, and Peter Battaglia. Polygen: An autoregressive
679 generative model of 3d meshes. In *International conference on machine learning*, pp. 7220–7229.
680 PMLR, 2020.
- 681 Songyou Peng, Chiyu Jiang, Yiyi Liao, Michael Niemeyer, Marc Pollefeys, and Andreas Geiger.
682 Shape as points: A differentiable poisson solver. *Advances in Neural Information Processing
683 Systems*, 34:13032–13044, 2021.
- 684 Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d
685 diffusion. In *The Eleventh International Conference on Learning Representations, ICLR 2023,
686 Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.net/
687 pdf?id=FjNys5c7VyY](https://openreview.net/pdf?id=FjNys5c7VyY).
- 688 Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets
689 for 3d classification and segmentation. In *CVPR 2017*, pp. 652–660, 2017a.
- 690 Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical fea-
691 ture learning on point sets in a metric space. *Advances in neural information processing systems*,
692 30, 2017b.
- 693 Scott Schaefer, Tao Ju, and Joe Warren. Manifold dual contouring. *IEEE Transactions on Visual-
694 ization and Computer Graphics*, 13(3):610–619, 2007.

- 702 Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra:
703 a hybrid representation for high-resolution 3d shape synthesis. In *Advances in Neural Information*
704 *Processing Systems (NeurIPS)*, 2021a.
- 705 Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra:
706 a hybrid representation for high-resolution 3d shape synthesis. *Advances in Neural Information*
707 *Processing Systems*, 34:6087–6101, 2021b.
- 708 Tianchang Shen, Jacob Munkberg, Jon Hasselgren, Kangxue Yin, Zian Wang, Wenzheng Chen, Zan
709 Gojic, Sanja Fidler, Nicholas Sharp, and Jun Gao. Flexible isosurface extraction for gradient-
710 based mesh optimization. *ACM Transactions on Graphics (TOG)*, 42(4):1–16, 2023.
- 711 Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view
712 diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023.
- 713 Yawar Siddiqui, Antonio Alliegro, Alexey Artemov, Tatiana Tommasi, Daniele Sirigatti, Vladislav
714 Rosov, Angela Dai, and Matthias Nießner. Meshgpt: Generating triangle meshes with decoder-
715 only transformers. *arXiv preprint arXiv:2311.15475*, 2023.
- 716 Jingxiang Sun, Bo Zhang, Ruizhi Shao, Lizhen Wang, Wen Liu, Zhenda Xie, and Yebin Liu.
717 Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior. *arXiv preprint*
718 *arXiv:2310.16818*, 2023.
- 719 Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative
720 gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023a.
- 721 Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm:
722 Large multi-view gaussian model for high-resolution 3d content creation. *arXiv preprint*
723 *arXiv:2402.05054*, 2024.
- 724 Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen.
725 Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. *arXiv preprint*
726 *arXiv:2303.14184*, 2023b.
- 727 Deemos Team. Deemos rodin. <https://hyperhuman.deemos.com/rodin>, 2024.
- 728 Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, Adam Letts, Yangguang Li, Ding
729 Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. Tripotr: Fast 3d object reconstruction
730 from a single image. *arXiv preprint arXiv:2403.02151*, 2024.
- 731 Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in*
732 *neural information processing systems*, 30, 2017.
- 733 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
734 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
735 *tion processing systems*, 30, 2017.
- 736 Peng-Shuai Wang, Yang Liu, and Xin Tong. Dual octree graph networks for learning adaptive
737 volumetric shape representations. *ACM Transactions on Graphics (TOG)*, 41(4):1–15, 2022.
- 738 Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolific-
739 dreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv*
740 *preprint arXiv:2305.16213*, 2023.
- 741 Zhengyi Wang, Yikai Wang, Yifei Chen, Chendong Xiang, Shuo Chen, Dajiang Yu, Chongxuan Li,
742 Hang Su, and Jun Zhu. Crm: Single image to 3d textured mesh with convolutional reconstruction
743 model. *arXiv preprint arXiv:2403.05034*, 2024.
- 744 Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan, Valentin Deschaintre, Kalyan Sunkavalli,
745 Hao Su, and Zexiang Xu. Meshlrn: Large reconstruction model for high-quality mesh. *arXiv*
746 *preprint arXiv:2404.12385*, 2024.
- 747 Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Sim-
748 plifying graph convolutional networks. In *International conference on machine learning*, pp.
749 6861–6871. PMLR, 2019.

- 756 Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and S Yu Philip. Multimodal large
757 language models: A survey. In *2023 IEEE International Conference on Big Data (BigData)*, pp.
758 2247–2256. IEEE, 2023.
- 759
760 Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. Instantmesh:
761 Efficient 3d mesh generation from a single image with sparse-view large reconstruction models.
762 *arXiv preprint arXiv:2404.07191*, 2024.
- 763
764 Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. Pointllm:
765 Empowering large language models to understand point clouds. *arXiv preprint arXiv:2308.16911*,
766 2023.
- 767
768 Xiaofeng Yang, Yiwen Chen, Cheng Chen, Chi Zhang, Yi Xu, Xulei Yang, Fayao Liu, and Guosheng
769 Lin. Learn to optimize denoising scores for 3d generation: A unified and improved diffusion prior
770 on nerf and 3d gaussian splatting. *arXiv preprint arXiv:2312.04820*, 2023.
- 771
772 Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from
773 one or few images. In *CVPR*, pp. 4578–4587, 2021.
- 774
775 Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. Sound-
776 stream: An end-to-end neural audio codec. *IEEE/ACM Transactions on Audio, Speech, and*
777 *Language Processing*, 30:495–507, 2021.
- 778
779 Longwen Zhang, Ziyu Wang, Qixuan Zhang, Qiwei Qiu, Anqi Pang, Haoran Jiang, Wei Yang, Lan
780 Xu, and Jingyi Yu. Clay: A controllable large-scale generative model for creating high-quality 3d
781 assets, 2024. URL <https://arxiv.org/abs/2406.13897>.
- 782
783 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo-
784 pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer
785 language models. *arXiv preprint arXiv:2205.01068*, 2022.
- 786
787 Zibo Zhao, Wen Liu, Xin Chen, Xianfang Zeng, Rui Wang, Pei Cheng, Bin Fu, Tao Chen, Gang Yu,
788 and Shenghua Gao. Michelangelo: Conditional 3d shape generation based on shape-image-text
789 aligned latent representation. *Advances in Neural Information Processing Systems*, 36, 2024.
- 790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A APPENDIX

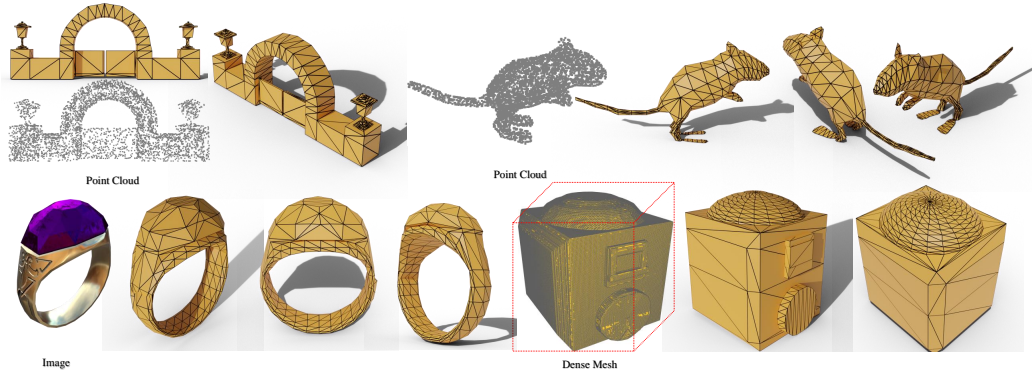


Figure 5: **Additional qualitative results of MeshAnything.** As shown, MeshAnything can be integrated with various 3D production pipelines to achieve highly controllable mesh generation.

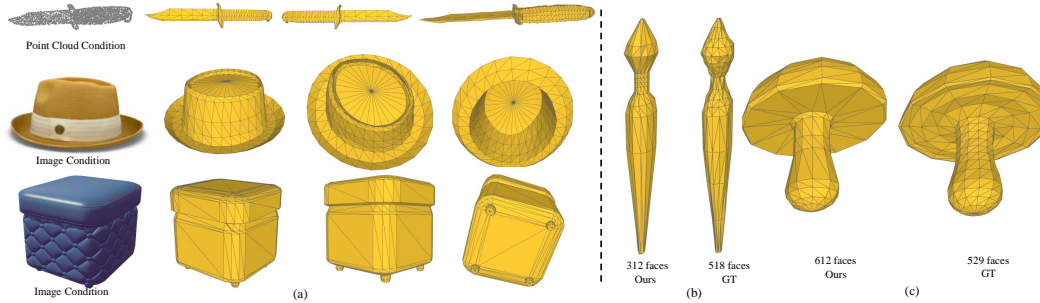


Figure 6: **Qualitative Results.** (a) further demonstrates our capability to achieve highly controllable mesh generation when combined with 3D asset production pipelines. Besides, we compare our results with ground truth in (b) and (c). In (b), MeshAnything generates meshes with better topology and 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.

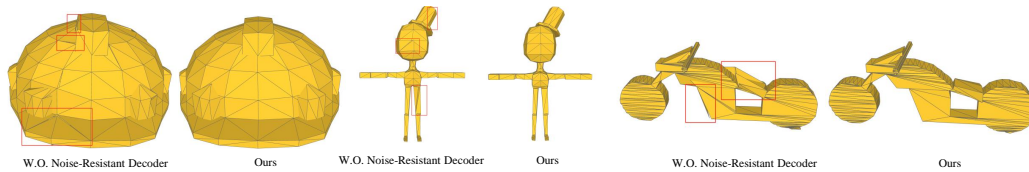


Figure 7: **Ablation on Noise-Resistant Decoder.** The decoder-only transformer may generate low-quality token sequences, and the decoder of VQ-VAE would typically produce flawed meshes based on these sequences. In contrast, our Noise-Resistant Decoder, aided by shape conditions, has the ability to resist these low-quality token sequences, producing higher-quality meshes.

A.1 METRICS

We follow the evaluation metric setting of (Siddiqui et al., 2023) in mesh generation experiments and the setting of (Chen et al., 2022) in mesh extraction experiments.

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.

Table 3: Reconstruction Performance under Different Noise Levels with and without Noise-Resistant (NR) Decoder. Please refer to A.1 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

Table 4: Ablation on Noise-Resistant (NR) Decoder for the Quality of Mesh Generation.

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

For mesh extraction, 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_R) and the same ratio for faces (#F_R).

For mesh generation, Coverage (COV) captures the diversity of generated meshes and is sensitive to mode dropping, but it does not reflect the quality of the results. Minimum Matching Distance (MMD) measures the average distance between the reference set and their nearest neighbors in the generated set, though it lacks sensitivity to low-quality outputs. The 1-Nearest Neighbor Accuracy (1-NNA) assesses both quality and diversity between the generated and reference sets. To evaluate topology quality, we render the ground truth meshes and generated meshes with their wireframes visualized. We then employ Frechet Inception Distance (FID) and Kernel Inception Distance (KID) on rendered images. MMD, and KID scores are scaled by a factor of 10^3 .

A.2 EXPERIMENTS

Additional Qualitative Experiments We present more qualitative results of MeshAnything here. As shown in Fig. 5 and Fig. 6, 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. 6 demonstrates that MeshAnything 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 (Lorenson & Cline, 1987) 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. 6 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.

Comparison with mesh extraction baselines. Our method is related to various mesh extraction methods (Lorenson & Cline, 1987; Chen & Zhang, 2021; Chen et al., 2022; Shen et al., 2023; Peng et al., 2021) since we also 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. Therefore, strictly speaking, our method cannot be considered the same as these reconstruction-based mesh extraction methods. The main purpose of this comparison

Table 5: Quantitative evaluation with mesh extraction baselines. MC, FC, SAP refer to Marching Cubes Lorensen & Cline (1987), FlexiCubes Shen et al. (2023), and Shape As Points Peng et al. (2021), respectively. Please refer to A.1 for metrics explanation.

Method	CD↓ ($\times 10^{-2}$)	ECD↓ ($\times 10^{-2}$)	NC↑	#V↓ ($\times 10^3$)	#F↓ ($\times 10^3$)	V.R↓	F.R↓
(a) Marching Cubes	1.532	6.733	0.954	73.22	146.0	440.2	462.2
(b) MC+Remesh (0.005)	2.174	7.813	0.912	127.8	167.9	748.1	534.6
(c) MC+Remesh (0.010)	2.083	7.578	0.929	39.01	41.78	225.4	132.3
(d) MC+Remesh (0.030)	2.915	8.329	0.863	5.848	4.410	34.38	14.05
(e) MC+Remesh (0.050)	4.179	8.138	0.814	2.299	1.538	13.64	4.920
(f) MC+Remesh (0.100)	7.312	10.771	0.748	0.625	0.359	3.735	1.149
(g) FC	1.190	6.121	0.967	59.12	121.1	378.2	391.1
(h) FC+Remesh (0.010)	1.861	6.940	0.933	37.98	40.19	205.5	124.2
(i) SAP	1.771	7.112	0.939	79.12	152.3	481.2	489.3
(j) SAP+Remesh (0.010)	2.367	7.862	0.925	39.17	42.87	239.1	136.6
(k) <i>MeshAnything</i>	2.256	6.245	0.902	0.172	0.318	0.888	0.871

is to use these mesh extraction methods as a reference for evaluating the quality of the meshes generated by MeshAnything in terms of shape. We compare MeshAnything with Lorensen & Cline (1987); Shen et al. (2023); Peng et al. (2021). Among these, MarchingCubes is the most popular mesh extraction method, FlexiCubes represents the state-of-the-art in mesh extraction, and Shape as Points is the leading method for extracting mesh from point cloud.

We also combined these methods with the remesh technique to test whether they could significantly reduce the number of faces while maintaining shape quality. We used Blender Remesh in voxel mode (Community, 2018; Blender Development Team, 2024), specifically using Blender version 4.1, as the remesh method. Since our evaluation dataset includes non-watertight meshes, we first extract the signed distance fields (SDF) of all ground truth meshes using (Wang et al., 2022), which can handle non-watertight meshes. We then apply Marching Cubes with a resolution of 128 on these SDFs. Next, we apply Blender remesh (Blender Development Team, 2024) 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. The settings of (Shen et al., 2023) and (Peng et al., 2021) follow their papers.

As shown in Tab. 5, we found that these methods require hundreds of times more faces to achieve results comparable to our method. Comparing (a), (g), (i) and (k), 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 these reconstruction-based mesh extraction methods. When comparing with remesh methods, we observe that they incur a high cost to achieve a face count similar to ours. Comparing (f) and (k), 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. It’s important to note that the metrics in mesh extraction can only indicate the quality of shape alignment, which do not effectively reflect the topological advantages of our method. 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.

Ablations on Noise-Resistant Conditional Decoder. We 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.

Table 6: Experiments on the Impact of Input Point Cloud Quality on Generated Results.

Method	CD↓ ($\times 10^{-2}$)	ECD↓ ($\times 10^{-2}$)	NC↑	#V↓ ($\times 10^3$)	#F↓ ($\times 10^3$)	V_R↓	F_R↓
(a) Noise scale 0.005	2.351	6.412	0.897	0.175	0.321	0.895	0.880
(b) Noise scale 0.020	2.980	6.970	0.881	0.180	0.330	0.901	0.910
(c) Noise scale 0.050	4.910	8.556	0.755	0.162	0.284	0.811	0.802
(d) Rodin	2.552	6.622	0.833	0.185	0.342	0.919	0.923
(e) <i>MeshAnything</i>	2.256	6.245	0.902	0.172	0.318	0.888	0.871

Table 7: **The comparison experiment between MeshAnything and MeshGPT on the chair category in ShapeNet.** This experiment is conducted with the same settings as in Table 2, except that both training and evaluation were performed exclusively on the chair category in ShapeNet.

Method	COV↑	MMD↓	1-NNA↓	FID↓	KID↓
MeshGPT	49.2	2.98	61.0	16.4	2.04
MeshAnything	62.1	1.92	49.8	11.2	1.21

After training both models for enough epochs, we test their performance to the same level of noise. As shown in Tab. 3, 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. 4, the model with Noise-Resistant Decoder achieved better results.

Experiments on the Impact of Input Point Cloud Quality on Generated Results. MeshAnything takes point clouds as input, and its robustness to point cloud quality determines its versatility across various applications. We design two experiments to evaluate its tolerance to input point cloud quality: First, keeping the other evaluation settings unchanged, we apply Gaussian noise to the input point cloud coordinates and normals. Specifically, for each point, we randomly sample Gaussian noise from a standard distribution, scale it by a noise factor, and add it to the point’s coordinates. The same approach is applied to the normals, but normalization is applied after adding the noise. Second, we use Rodin’s generation result as the ground truth mesh, sample point clouds from this mesh as input, and evaluate the deviation between the generated result and the ground truth.

As shown in Tab. 6, MeshAnything did not experience a significant performance drop in (a) and (b), demonstrating resilience to noise in the point cloud, with a noticeable performance decrease only in (c). It is important to note that the input point cloud is normalized to the range $[-1, 1]$, and the noise scale in (c) is already quite large. The experiment in (d) further demonstrates that MeshAnything can tolerate generated point clouds and effectively integrate with 3D generation models.

Experiments between MeshAnything and MeshGPT on the chair category in ShapeNet. To reduce the complexity of learning the shape distribution, we conducted the comparison experiment between MeshAnything and MeshGPT (Siddiqui et al., 2023) exclusively on the chair category in ShapeNet. As shown in 7, the gap between MeshGPT and MeshAnything narrowed under this setting, though a noticeable difference remains. This result supports our argument in Sec. 3 that shape-conditioned mesh generation is considerably less challenging to learn than unconditional mesh generation.

Sampling Strategy and Diversity of Generated Results All of our qualitative and quantitative experiments employed multinomial sampling, which is equivalent to a beam search with a beam size of 1. We used top-k and top-p values set to 50 and 0.95, respectively. We provide the generation variants obtained using this sampling approach in Fig. 8. As shown in the figure, MeshAnything demonstrates the capability to generate various topologies.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

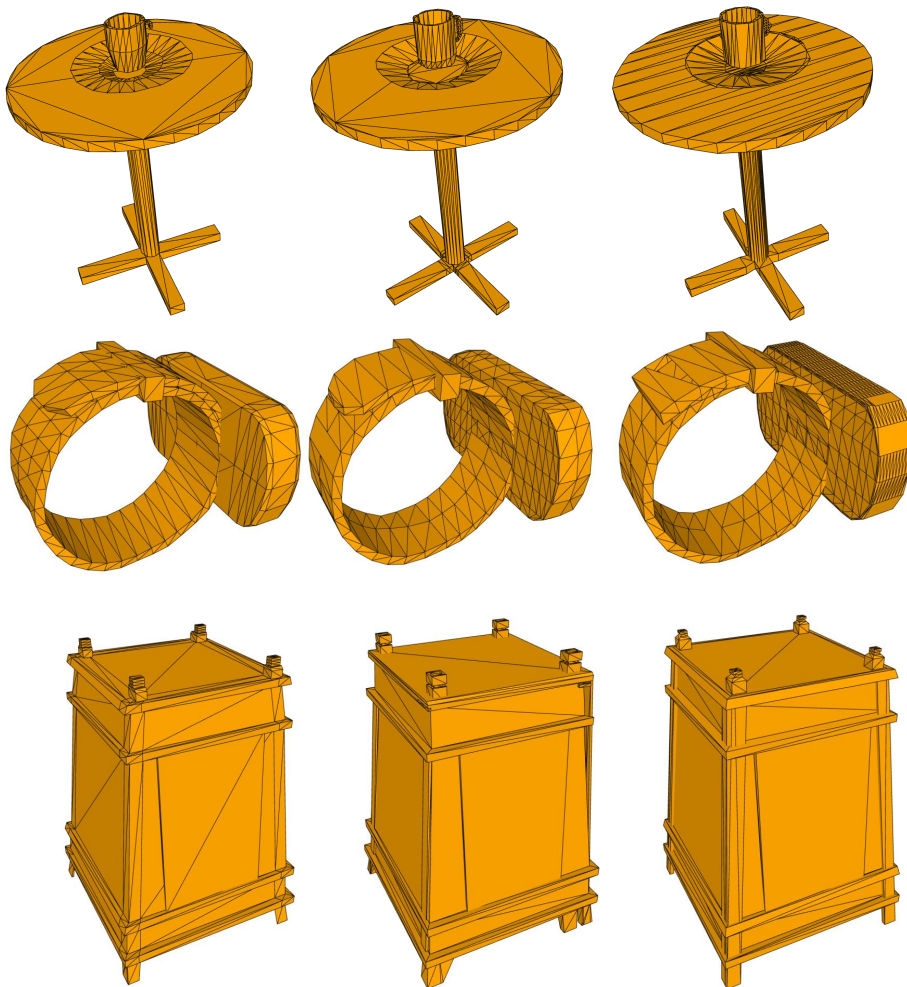


Figure 8: **Diversity of Generated Results.** The generated results of MeshAnything obtained using multinomial sampling with different random seeds.

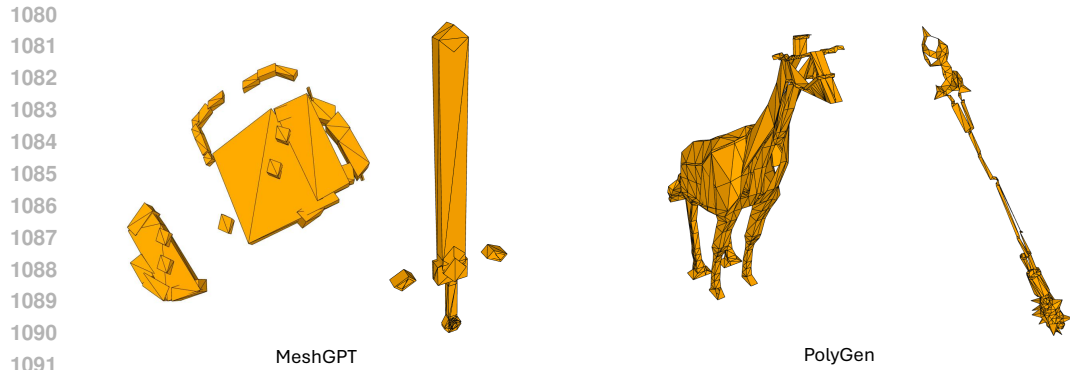
Qualitative Comparison with MeshGPT. We provide a quantitative comparison with MeshGPT (Siddiqui et al., 2023) and Polygen (Nash et al., 2020). As shown in Fig. 9, it is challenging for these two methods to produce high-quality meshes on Objaverse. This advantage is primarily because MeshAnything avoids the need to learn the complex 3D shape distribution, making the training process significantly easier compared to unconditional mesh generation.

A.3 LIMITATIONS

Our method cannot generate meshes that exceed the maximum face count limit, which restricts its ability to convert large scenes and highly complex objects into meshes. Additionally, due to its generative nature, our method is not as stable as reconstruction-based mesh extraction methods like (Lorenzen & Cline, 1987; Shen et al., 2023). After decades of development, reconstruction methods achieve near 100% success rates with reasonable inputs. However, as a generative approach, our method inevitably produces occasional failure cases.

A.4 FUTURE WORKS.

Autoregressive mesh generation remains inefficient, significantly limiting the practical application of this approach. Developing more efficient representations for mesh generation is essential. Fur-



1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Figure 9: **Qualitative Comparison with MeshGPT (Siddiqui et al., 2023) and PolyGen (Nash et al., 2020)**. Since MeshAnything is conditioned on point cloud input, a direct qualitative comparison with these two methods is not feasible. Therefore, we showcase unconditional generation results from MeshGPT and PolyGen. For a fair comparison, both methods use the same architecture and data as our approach. As shown, when scaled to Objaverse, MeshGPT and PolyGen struggle to consistently produce meshes with complex, general shapes.

thermore, improving the success rate of mesh generation is critical. Due to the cumulative error characteristic of autoregressive models, a deviation in any one token prediction can cause the entire mesh generation process to fail. Future developments in this line of research may benefit from techniques used in large language models (LLMs) to mitigate similar issues. Additionally, current mesh generation methods are limited to the object level. Extending these methods to the scene level is equally important, as it would greatly expand the range of potential applications.

A.5 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.