

## Step1X-3D: Towards High-Fidelity and Controllable Generation of Textured 3D Assets

**Step1X-3D Team & LightIllusions Team**  
StepFun & LightIllusions  
<https://github.com/stepfun-ai/Step1X-3D>

## Abstract

While generative artificial intelligence has advanced significantly across text, image, audio, and video domains, 3D generation remains comparatively underdeveloped due to fundamental challenges such as data scarcity, algorithmic limitations, and ecosystem fragmentation. To this end, we present Step1X-3D, an open framework addressing these challenges through: (1) a rigorous data curation pipeline processing >5M assets to create a 2M high-quality dataset with standardized geometric and textural properties; (2) a two-stage 3D-native architecture combining a hybrid VAE-DiT geometry generator with an diffusion-based texture synthesis module; and (3) the full open-source release of models, training code, and adaptation modules. For geometry generation, the hybrid VAE-DiT component produces TSDF representations by employing perceiver-based latent encoding with sharp edge sampling for detail preservation. The diffusion-based texture synthesis module then ensures cross-view consistency through geometric conditioning and latent-space synchronization. Benchmark results demonstrate state-of-the-art performance that exceeds existing open-source methods, while also achieving competitive quality with proprietary solutions. Notably, the framework uniquely bridges the 2D and 3D generation paradigms by supporting direct transfer of 2D control techniques (e.g., LoRA) to 3D synthesis. By simultaneously advancing data quality, algorithmic fidelity, and reproducibility, Step1X-3D aims to establish new standards for open research in controllable 3D asset generation.



Figure 1: Step1X-3D demonstrates the capability to generate 3D assets with high-fidelity geometry and versatile texture maps, while maintaining exceptional alignment between surface geometry and texture mapping. From left to right, we sequentially present: the base geometry (untextured), followed by cartoon-style, sketch-style, and photorealistic 3D asset generation results.

# 1 Introduction

In recent years, generative artificial intelligence (GAI) technology has achieved remarkable development across five primary domains: text [1, 73, 42, 35], image [23, 62, 60, 33, 28, 22, 44], audio [76, 26, 19], video [46, 51, 75], and 3D content [54, 32, 31, 55, 24, 92]. However, unlike these other modalities, 3D generation technology exhibits significantly lower maturity and slower development progress, remaining far from production-ready. To facilitate further advancements in the field of 3D generation, we analyze the main challenges of 3D content generation from the respective datasets, algorithms, and their corresponding technological ecosystems.

Firstly, data scarcity constitutes the primary bottleneck in 3D generation advancement, where utility is governed by both quantity and quality. Current open-source datasets with >10K samples are limited to ShapeNet [6] (51K), Objaverse [18] (800K), and Objaverse-XL [17] (10.2M). Although Objaverse-XL exceeds 10M samples, its web-sourced provenance results in unavoidable quality heterogeneity. Secondly, the inherent complexity of 3D representations—where geometry and texture are decoupled—renders quality assessment fundamentally more challenging than in other modalities. Finally, despite rapid advances, the 3D generation ecosystem remains underdeveloped, with a growing gap between opensource [38, 85, 94] and proprietary solutions. For instance, an open-source model like Trellis [85] suffers from limited generalization due to its 500K-scale training dataset. Meanwhile, the release strategy for some advanced models, like Hunyuan3D 2.0 [71], providing only pre-trained weights (without training code), can restrict fine-tuning. Furthermore, these models often lack the conditional generation support typically found in commercial platforms like Tripo and Rodin. These challenges in data availability, reproducibility, and controllability significantly impede progress in 3D generation technologies.

To address this disparity, we introduce Step1X-3D, a 3D native framework that combines advanced data curation with an innovative two-stage architecture as shown in Fig. 2. Our goal is to achieve high-fidelity, controllable 3D asset generation while championing reproducibility through the open release of curated data and training methodologies. The Step1X-3D data pipeline begins with over 5M 3D assets sourced from public datasets (e.g., Objaverse [18], Objaverse-XL [17], and etc) and proprietary collections. These assets undergo a rigorous multi-stage refinement process: first, low-quality textures are eliminated based on criteria like resolution, normal map accuracy, material transparency, and surface complexity; second, watertight mesh conversions are enforced to ensure geometric consistency for robust training supervision. This process yields a curated dataset of 2M high-quality assets, a subset of which (almost 800K assets derived from public data) will be openly released. Architecturally, Step1X-3D employs: (1) a geometry generation stage using a hybrid 3D VAE-DiT diffusion model to produce Truncated Signed Distance Function (TSDF) (later meshed via marching cubes), and (2) a texture synthesis stage leveraging an SD-XL-fine-tuned multi-view generator. This generator conditions on the produced geometry and input images to produce view-consistent textures that are then baked onto the mesh. This integrated approach aims to advance the field by simultaneously resolving critical challenges in data quality, geometric precision, and texture fidelity, while establishing new benchmarks for open, reproducible research in 3D generation.

The Step1X-3D geometry generation framework employs a latent vector set representation [91] to encode point clouds into a compact latent space, which is decoded into TSDF through a scalable perceiver-based encoder-decoder [29, 95] architecture. To preserve high-frequency geometric details, we incorporate sharp edge sampling and integrate dual cross attention mechanisms derived from DoRA [9]. For the diffusion backbone, we adapt the state-of-the-art MMDiT architecture from FLUX [33] – originally developed for text-to-image generation – by modifying its transformer layers to process our 1D latent space. This VAE-Diffusion hybrid design, architecturally analogous to contemporary 2D generative systems, facilitates direct transfer of 2D parameter-efficient adaptation methods (e.g., LoRA [88]) to 3D mesh synthesis. Consequently, our framework uniquely enables single-view-conditioned pretraining on large-scale datasets while maintaining compatibility with established 2D fine-tuning paradigms for downstream 3D generation tasks, effectively bridging 2D and 3D generative approaches.

The Step1X-3D texture synthesis pipeline initiates with post-processing of the Step1X-3D geometry output using Trimesh to rectify surface artifacts (including non-watertight meshes, topological irregularities, and surface discontinuities), followed by UV parameterization via xAtlas [52]. The synthesis process employs a three-stage architecture: (1) a multi-view image generation diffusion model that conditions on both input images and rendered geometric maps (normal and position)

to enforce view consistency and geometric alignment; (2) a texture-space synchronization module integrated within the denoising process to maintain cross-view coherence through latent space alignment; and (3) texture completion via multi-view back-projection with subsequent texture-space inpainting to address occlusion artifacts and generate seamless UV maps. This hierarchical approach ensures both geometric fidelity and photometric consistency throughout the texture generation pipeline.

In summary, the key contributions of this technical report are threefold:

- We present a comprehensive data curation pipeline that provides insights into 3D asset characteristics while enhancing generation fidelity.
- We propose Step1X-3D, a 3D native generation framework that decouples geometry and texture synthesis. It produces topologically-sound meshes with geometrically-aligned textures, while enabling enhanced controllability through image and semantic inputs. The full framework - including base models, training code, and LoRA-based adaptation modules - will be open-sourced to benefit the 3D research community.
- Extensive comparative experiments demonstrate that Step1x-3D surpasses existing open-source 3D generation approaches in asset quality while achieving performance comparable to proprietary state-of-the-art solutions.

## 2 Related work

### 2.1 Optimization-based 3D Generation

Optimization-based 3D generation is a framework where 3D representations and their optimization procedures are independently defined across distinct objects, enabling text-to-3D synthesis through input prompts. This approach incorporates diverse 3D representations, ranging from geometric primitives (e.g., voxels, meshes [4]) to advanced neural representations such as Neural Radiance Fields (NeRF) [54] and 3D Gaussian Splatting (3D GS) [32]. Text prompts are encoded into pre-trained models (e.g., CLIP [59], Stable Diffusion [62] and etc) to extract semantic-aligned 2D visual features, which supervise the multi-view consistency of rendered 3D representations via differentiable rendering pipelines. By iteratively optimizing geometric and appearance parameters through gradient-based alignment of 2D supervision, the framework achieves cross-modal (text-to-3D) consistence and produces 3D assets. These optimization-based 3D generation methods [30, 64, 58, 8, 97], while capable of processing arbitrary text prompts, suffer from limited support for image-conditioned inputs and inherent limitations such as prolonged optimization cycles, poor generalization of trained parameters and inferior generation quality. These shortcomings critically hinder their practical deployment and scalability in real-world applications.

### 2.2 Feed-forward 3D Generation

To enhance generation efficiency and quality, feed-forward 3D generation methods have emerged as a promising alternative. Specifically, feed-forward methods aim to fit 3D datasets comprising 3D representations (e.g., multi-view images, point cloud, NeRF, 3D GS and mesh), enabling direct 3D synthesis via forward propagation during inference without iterative optimization, thereby achieving order-of-magnitude acceleration compared to optimization-based approaches. Concurrently, advancements in 3D data availability and algorithmic innovations have driven qualitative leaps in generation fidelity and generalization capability. The evolution of feed-forward 3D generation can be categorized into two distinct paradigms: *2D-lifting-to-3D generation paradigm* and *3D native generation paradigm*.

**2D-lifting-to-3D Generation Paradigm** The 2D-lifting-to-3D paradigm introduces 2D priors into feedforward frameworks, analogous to optimization-based methods, allowing cross-modal knowledge transfer from large-scale 2D datasets. The 2D-lifting-to-3D paradigm is broadly divided into two categories, with the first focusing on single-stage 3D reconstruction from single image input through image embedding techniques (e.g., DINOv2 [56]) for feature extraction, followed by 3D geometry synthesis via Vision Transformer (ViT)[20] architectures. Representative methods include LRM [24], which reconstructs Neural Radiance Fields (NeRF) [54] from encoded image features, TGS [98] that

generates 3D Gaussian Splatting (3D GS) [32] representations via transformer-based decoders, and TripoSR [72] and MeshLRM [81], which leverage hybrid decoders with differentiable rasterization to directly output mesh structures, collectively demonstrating the versatility of this paradigm across diverse 3D representations. The second category adopts a two-stage 2D-lifting-to-3D paradigm, explicitly prioritizing multi-view consistency for robust 3D generation. These methods accept text prompts or single-image inputs and leverage pre-trained 2D stable diffusion models [62] to synthesize multi-view images, which are subsequently fed into reconstruction networks to produce diverse 3D representations. Research efforts focus on dual optimization: enhancing multi-view synthesis quality in the first stage and improving 3D reconstruction fidelity in the second. For multi-view consistency, SyncDreamer [47] introduces cost volume constraints to align geometry across views, Wonder3D [49] jointly generates RGB images and normal maps for cross-modal consistency, Era3D [36] employs epipolar attention to boost computational efficiency while scaling output resolution from 256×256 to 512×512, and MV-Adapter [27] further extends resolution to 768×768. Beyond 2D image diffusion priors, SV3D [74] and VideoMV [99] leverage stable video diffusion models [3] with temporal coherence to achieve superior multi-view alignment. Reconstruction-oriented methods like LGM [66] and GRM [87] refine neural field representations, while InstantMesh [86], CRM [79], and Unique3D [83] prioritize practical mesh reconstruction through differentiable rasterization. Approaches including MV-Diffusion++ [68], Meta 3D Gen [2], and Cycle3D [69] separately enhance algorithmic components across both stages—refining multi-view synthesis pipelines for consistency and optimizing reconstruction networks for geometric fidelity—through meticulous design of each standalone module within the two-phase framework. Since 2d-lifting-to-3d generation methods mainly rely on 2D images as their foundational supervision, they predominantly prioritize the perceptual quality of generated images while neglecting 3D geometric fidelity. Consequently, the resulting 3D geometry frequently exhibits incompleteness and a lack of fine-grained detail.

**3D Native Generation Paradigm** Through explicit geometric representation modeling and learned 3D feature extraction, native 3D generation frameworks enable precise geometry synthesis. In the primitive 3D era, early works like Point-E [55] and Shape-E [31], constrained by limited 3D datasets [5] and underdeveloped architectures, were restricted to limited-category object generation with low visual fidelity. Current native 3D models primarily follow the CG modeling pipeline, decomposing 3D generation into a two-stage process: first generating geometry, then using the geometry to guide texture generation. These native 3D models directly model the geometric component of 3D representations through generative models, achieving more precise and accurate geometric results compared to 2D-lifting-to-3D approaches [24, 66, 79, 49, 86, 70]. These significant advancements are largely attributed to recent progress in 3D data and algorithmic developments. On the data side, early native 3D algorithms were primarily trained on ShapeNetCore [5], which contained only about 51,300 objects. Later, with the release of Objaverse [18] and Objaverse-XL [18], the dataset sizes expanded to 800K+ and 10M+, respectively, greatly enriching the 3D asset database. This played a crucial role in validating model generation capabilities and scalability. Many proprietary algorithms (e.g., Tripo, Rodin and Hunyuan) have even incorporated additional private data to further enhance generalization or specialized performance. Beyond data, algorithmic advancements have also been pivotal. In 3D Variational Auto-Encoder (VAE), Michelangelo [95] innovatively introduced a perceiver-based [29] architecture for 3D point cloud feature extraction while aligning 3D with text and image modalities. Direct3D [84] encodes 3D shapes into a latent triplane space. 3DShape2VecSet [91] proposed a transformer-based method for representing 3D latent space. CLAY [92] further scale 3DShape2VecSet to a large scale dataset and demonstrate the unprecedented potential of 3D native diffusion model. Trellis [85] introduced a Structured Latent Representation by aggregating features from 3D voxels and multi-view inputs. Dora [9] proposed a sharp edge sampling strategy to improve geometric detail reconstruction. In generative algorithms, 3D diffusion models have been constructed by adapting 2D generative techniques such as flow matching [41, 45] and DiT [57] architectures, ensuring robust conditional injection and geometric accuracy. Meanwhile, autoregressive 3D generation methods have also seen rapid development, with some focusing purely on generation (e.g., Mesh-GPT [65], Mesh-XL [10], LLAMA-mesh [78] and the recent OctGPT [80]), while others (e.g., MeshAnything [11], EdgeRunner [67] and BPT [82]) optimize dense meshes for intelligent retopology to facilitate downstream tasks like shading and rendering. Currently, mainstream 3D asset generation tools—such as the open-source CraftsMan3D [38], Trellis [85] and Hunyuan3D 2.0 [71], and proprietary solutions like Tripo, Rodin [92], and Meshy—primarily rely on 3D diffusion for geometry generation, as autoregressive methods still lag in effectiveness due to limitations in token representation. Beyond geometry generation, texture synthesis has also made

significant strides. TEXTure [61] and Text2Tex [7] leverage pre-trained depth-to-image diffusion models [93], iteratively generating multi-view images from given camera trajectories and depth maps for texture baking. However, independently generated views suffer from severe consistency issues. Paint-3D [90] adopts a coarse-to-fine approach, introducing UV inpainting and UVHD diffusion models to refine incomplete regions and remove illumination artifacts in UV space. SyncMVD [48] enhances multi-view consistency by incorporating a texture synchronization module in latent space, while MVPaint [12] argues that latent space resolution ( $32 \times 32$ ) is insufficient for fine texture details and instead performs texture synchronization in image space with resolution ( $128 \times 128$ ), combined with 3D point cloud inpainting. This technique report follows the 3D diffusion generation paradigm to generate high-fidelity 3D shapes and textures compared to other state-of-the-art (SOTA) methods, and further involve controllable modules to improve the flexibility of the generation process.

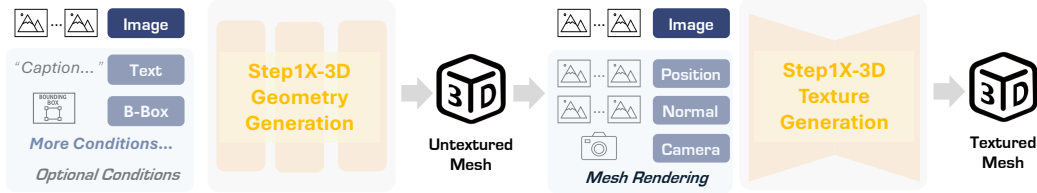


Figure 2: Overall pipeline of Step1X-3D.

### 3 Step1X-3D Geometry Generation

Step1X-3D is a flow-based diffusion model designed to generate high-fidelity 3D shapes from images, with support for multi-modal conditioning including text and semantic labels. The proposed geometry generation model builds upon prior latent set diffusion models such as Shape2VecSet [91], CLAY [92], Michelangelo [95], and Craftsman3D [38], utilizing a latent set diffusion framework with rectified flow for 3D shape generation.

In this section, we first introduce the data curation methodology for preprocessing in Sec. 3.1. Next, we provide a detailed description of the architectural design for both the Shape VAE and diffusion model components in Sec. 3.2. Additionally, inspired by the approach in CLAY [92], we present an adaptation of the LoRA [25] ecosystem for 3D generation in Sec. 3.3. All training code and sampled data will be made publicly available to support research and community development.

#### 3.1 Geometry Data Curation

Recent years have witnessed the emergence of several large-scale open-source 3D datasets, including Objaverse [18], Objaverse-XL [17], ABO [14], 3D-FUTURE [21], ShapeNet [5], etc., which together contain more than 10 million 3D assets. However, as most of this data is sourced from the web—particularly the extensive Objaverse-XL collection—the quality varies considerably. To ensure the data is suitable for training, we implemented a comprehensive 3D data processing pipeline that performs thorough preprocessing to curate a high-quality, large-scale training dataset.

Our curation pipeline consists of three main stages. First, we filter out low-quality data by removing assets with poor textures, incorrect normals, transparent materials or single surface. Second, we convert non-watertight meshes into watertight representations to enable proper geometry supervision. Third, we uniformly sample points on the surface along with their normals to provide comprehensive coverage for the VAE and diffusion model training. Through our comprehensive data processing pipeline, we successfully curated roughly 2 million high-quality 3D assets from multiple sources: extracting 320k valid samples from the original Objaverse dataset, obtaining an additional 480k from Objaverse-XL, and combining these with carefully selected data from ABO, 3D-FUTURE, and some internal created data.

**Data Filter** The complete data filter process is shown in Fig. 3 (a). (1) *Texture Quality Filtering*: We render 6 canonical-view albedo maps for each 3D model. These rendered images are then converted to HSV color space for analysis. For each view, we compute histograms of the Hue (H) and Value (V) channels. Based on these histograms, we filter out textures that are either too dark,

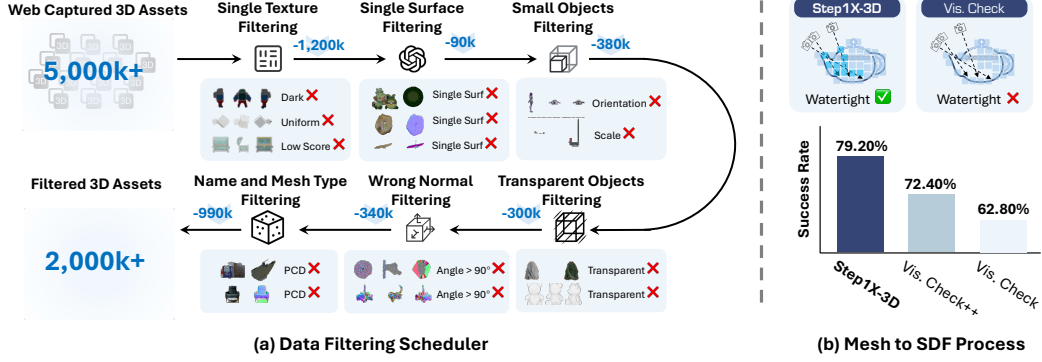


Figure 3: **Data curation pipeline.** (a) we filter bad quality data from web captured 3D assets (e.g., Objaverse, ABO, Thingiverse, and self-collected assets). (b) improved conversion rate of mesh-to-SDF.

too bright, or overly uniform in color. We then compute perceptual scores for these six views and sort the data accordingly, removing the bottom 20% of lowest-ranked samples. (2) *Single-surface Filtering*: We render 6 canonical-view Canonical Coordinate Maps (CCM) [77, 37] to detect single-surface geometry. Specifically, we check whether corresponding pixels on opposite views map to the same 3D point. If the ratio of such pixel matches exceeds a threshold, the object is classified as single-surfaced. (3) *Small Object Filtering*: We filter out data where the target object occupies too small an area in frontal views. This occurs in two scenarios: improper object orientation (e.g., a supine human where only feet are visible in front view), or distant objects in multi-object scenes that become too small after normalization. Specifically, we calculate the percentage of valid alpha-channel pixels in frontal views and discard samples with less than 10% pixels coverage. (4) *Transparent Object Filtering*: We exclude objects with transparent materials, as they are typically modeled using alpha-channel planes (e.g., for tree leaves). These transparent surfaces cause misalignment between rendered RGB images and actual geometry, adversely affecting model training. Our filtering method detects and removes assets whose Principled BSDF shaders contain alpha channels. (5) *Wrong Normal Filtering*: We identify and remove data with incorrect normals, which would otherwise create holes during watertight conversion. Our method renders 6-view normal maps in camera space and detects erroneous normals by checking whether any normal vectors form obtuse angles with their corresponding camera position. (6) *Name and Mesh Type Filtering*: We also filter out data labeled as point clouds by name or mesh type, as these scan-derived datasets typically contain noisy geometry and are difficult to convert into watertight meshes.

**Enhanced mesh-to-SDF** Training a Shape VAE [91, 9] requires watertight mesh to enable the extraction of SDF (Signed Distance Function) field from the processed meshes, which serve as geometric supervision [92]. CLAY [92] introduced a "visibility check" method for Mesh-to-SDF conversion, which first splits the space into  $N^3$  grids. Each grid center checks the visibility using depth test and utilizes a mask with size  $N^3$  to indicate whether the grid is considered invisible. However, for the non-manifold objects with holes, like the windows on the wall, it is easy to encounter floaters inside the converted mesh. To address this challenge, we implement a robust classification scheme by incorporating the concept of the winding number [53] as in CraftsMan3D [38], which is an effective tool for determining whether points are inside or outside a shape. For each point sampled within the voxel grid, we compute its generalized winding number<sup>1</sup>, considering points with values above our empirically determined threshold of 0.75. The resulting winding number mask is then combined with the original visibility test through logical conjunction to generate the final occupancy mask for MarchingCubes [50]. Experimental results as show in Fig. 3 (b) demonstrate this approach achieves a 20% watertight conversion success rate improvement on the Objaverse dataset [18].

**Training Data Conversion** (1) Data for VAE: Following the Dora [9], we employ the Sharp Edge Sampling (SES) strategy to enhance point sampling in geometrically salient regions. Specifically, we combine uniformly sampled points  $P_{uniform}$  with additional points  $P_a$  sampled from salient areas,

<sup>1</sup><https://libigl.github.io/tutorial/#generalized-winding-number>

forming the final point set  $P = P_{uniform} \cup P_{salient}$  along with their corresponding normals, as input to the VAE. For geometry supervision, we sample three distinct sets of points with their SDF values: 200k points within the cube volume, 200k points near the mesh surface with a threshold 0.02, and 200k points directly on the surface. (2) Data for Diffusion: For training our single-image conditioned flow model, we render each 3D model from 20 randomly sampled viewpoints with camera elevation between  $-15^\circ$  and  $30^\circ$ , azimuth between  $-75^\circ$  and  $75^\circ$ , and focal lengths randomly selected from orthogonal projection or perspective projections (with focal length uniformly sampled from 35mm to 100mm). We adjust the camera position to ensure the content occupies approximately 90% of the image. Additionally, we apply common data augmentations such as random flipping (for both images and sampled meshes), color jitter, and random rotations between  $-10^\circ$  and  $10^\circ$ .

### 3.2 Step1X-3D Shape Generation

Similar to 2D image generation [63, 33], the Step1X-3D shape generation module consists of a shape autoencoder and a Rectified Flow Transformer. For the sampled point cloud  $P$ , we first compress it into a 1D tensor using a shape latent set autoencoder [91], then train the diffusion model with a 1D Rectified Flow Transformer inspired by Flux [33]. We also support additional components like LoRA for extended flexibility.

**3D Shape Variational Autoencoder** The success of the Latent Diffusion Model (LDM) [63] proves that a compact, efficient, and expressive representation is essential for training a diffusion model. Therefore, we first encode 3D shapes into a latent space and then train a 3D latent diffusion model for 3D generation. Following the design of 3DShape2VecSet [91], we adopt a latent vector set representation to encode point clouds into latent space and decode them into geometric functions (e.g., signed distance fields or occupancies). To improve scalability, we adopt a transformer-based encoder-decoder architecture as in recent works [29, 95]. Additionally, we incorporate the Sharp Edge Sampling and Dual Cross Attention techniques proposed in Dora [9] to enhance geometric detail preservation. Specifically, we use the downsampled variant of 3DShape2VecSet. Instead of learnable queries, we initialize the latent queries  $S = \text{FPS}(P_{uniform}) \cup \text{FPS}(P_{salient})$  directly with the point cloud itself using the Farthest Points Sampling (FPS). We first integrate the information of the concatenated Fourier positional encodings with their respective normals into the shape encoder, forming the actual input to the shape encoder:  $\hat{P} = \text{Concat}(PE(P_c), P_n)$ , where  $P_c$  is the points position and  $P_n$  is the normal. The encoder then processes this input using two cross-attention layers and  $L_e$  self-attention layers, encoding the points into the latent space via:

$$\begin{aligned} \text{Enc}(P) &= \text{SelfAttn}^{(i)}(\text{CrossAttn}(S, \hat{P}_{uniform}), \text{CrossAttn}(S, \hat{P}_{salient})), \\ &\forall i = 1, 2, \dots, L_e. \end{aligned} \quad (1)$$

Similarly, we use a perceiver-based decoder that mirrors the architecture of the encoder, and has an extra linear layer  $\varphi_O$  to learn to predict the Truncated Signed Distance Function (TSDF) value at  $x$ :

$$\begin{aligned} \text{Dec}(x|S) &= \varphi_O(\text{CrossAttn}(PE(x), \text{SelfAttn}^{(i)}(S))), \\ &\forall i = 1, 2, \dots, L_d, \end{aligned} \quad (2)$$

where  $L_d$  is the number of self-attention layers in the shape decoder. Given a query point  $x \in \mathbb{R}^3$  in 3D space and a learned latent set  $S$ , the decoder can output its TSDF value. Then the training objective is as:

$$\mathcal{L}_{VAE} = \mathbb{E}_{x \in \mathbb{R}^3} \left[ \text{MSE} \left( \hat{O}(x|S), \text{Dec}(x) \right) \right] + \lambda_{kl} \mathcal{L}_{kl}, \quad (3)$$

where  $\hat{O}(x)$  is the ground truth TSDF value of  $x$  and the truncated scale is set to  $2/256$ . The KL divergence loss  $\mathcal{L}_{kl}$  is used to regularize the latent space distribution to a standard Gaussian distribution. Subsequently, we sample query points from a regular grid to obtain their corresponding TSDF values, which are then utilized to reconstruct the final surface using the Marching Cubes [50]. We also employ Hierarchical Volume Decoding [34] to accelerate the inference process. Please refer to the 3DShape2VecSet [91] and Dora [9] for more details.

**Step1X-3D Diffusion Backbone** Following the state-of-the-art text-to-image Diffusion model architectures FLUX [33], we use the same MMDiT structure, but modify it for 1D latent space

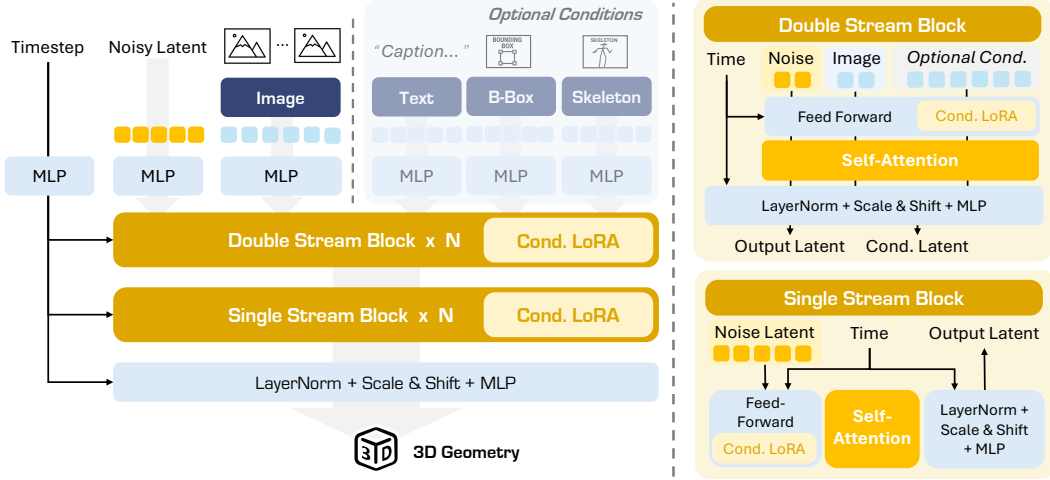


Figure 4: **Step1X-3D geometry diffusion framework.** The framework enables seamless integration of diverse conditions, allowing fine-tuning for various downstream tasks.

processing as illustrated in Fig. 4. In the dual-stream blocks, latent tokens and condition tokens are processed separately with their own QKV projections and MLPs, but they still interact through cross-attention. In contrast, single-stream blocks combine both types of tokens and process them jointly using parallel spatial and channel attention mechanisms. This hybrid approach allows flexible feature learning while maintaining efficient cross-modal interaction. Notably, to effectively introduce spatial position information in different noise patches, FLUX.1 employs rotary position encoding (RoPE) to encode spatial information within each noise patch. Since our ShapeVAE’s latent sets representation lacks explicit spatial correspondence, we remove positional embeddings for the latent sets  $S$  and retain only timestep embeddings for modulation purposes. For single-image conditioned shape generation, we leverage a pre-trained DINOv2 large image encoder [56] with registers [15] to extract conditional tokens from preprocessed 518×518 resolution images - where we perform background removal, object centering/resizing, and white background filling to enhance effective resolution and minimize background interference. To capture both semantic and global image characteristics, we concatenate complementary features from CLIP-ViT-L/14 [64]. These combined features are then injected through parallel cross-attention mechanisms within each flow block, enabling simultaneous processing of both global and local visual information.

### 3.3 More flexible control for 3D generation

Building upon the structural advantages of our VAE with Diffusion framework - which mirrors contemporary text-to-image architectures - we achieve seamless transfer of 2D controllable generation techniques (e.g., ControlNet, IP-Adapter) and parameter-efficient adaptation methods like LoRA to 3D mesh synthesis. As demonstrated by CLAY [92] in exploring ControlNet-UNet combinations for 3D generation, we systematically implement these control mechanisms within our Step1x-3D framework and support diffusers community<sup>2</sup>. To efficiently incorporate conditional signals while preserving the pre-trained model’s generalization ability, we can introduce an additional Condition Branch using the ControlNet-like strategy or LoRA. During the current open-source phase, we implement LoRA for geometric shape control using label as reference examples. This lightweight solution achieves domain-specific fine-tuning with minimal resource overhead while preserving the original feature space integrity, resulting in controllable shape generation. We first annotate each mesh with geometric attributes, such as symmetry or level of geometric detail. Using these annotations, we can train an additional LoRA module rather than fine-tuning the entire network, enabling the model to utilize these labels for controlling object geometry. This LoRA adaptation exclusively to the Condition Branch, enabling efficient injection of conditional signals without compromising the pretrained model’s capabilities. We will also plan introduce updates in later stages to incorporate fine-tuning through skeleton, bounding box (bbox), caption, and IP-image conditions.

<sup>2</sup><https://huggingface.co/docs/diffusers/index>

### 3.4 Training Rectified Flow Model

For training, we utilize a flow matching objective that constructs a probability path between Gaussian noise  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and data distributions, where Rectified Flow’s linear sampling mechanism streamlines network training by directly predicting the velocity field  $\mathbf{u}_t = \frac{d\mathbf{x}_t}{dt}$  that transports samples  $\mathbf{x}_t$  toward target data  $\mathbf{x}_1$ , thereby improving both efficiency and training stability. Building on SD3’s logit-normal sampling strategy, we strategically increase the sampling weight for intermediate timesteps  $t \in (0, 1)$  during training, as these mid-range temporal coordinates pose greater prediction challenges for velocity estimation in the Rectified Flow framework. The final objective is formulated as:

$$\mathcal{L} = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}_1} [\|\mathbf{u}_\theta(\mathbf{x}_t, \mathbf{c}, t) - \mathbf{u}_t\|_2^2] \quad (4)$$

where  $\mathbf{c}$  denotes the conditioning signal, with an adaptive time-step weighting scheme. To further stabilize training, we incorporate an Exponential Moving Average (EMA) strategy with a decay rate of 0.999 to smooth parameter updates. The training is conducted in two phases: initially, for rapid convergence, we use a latent set size of 512, a learning rate of 1e-4, and a batch size of 1920 across 96 NVIDIA A800 GPUs for 100k iterations. Subsequently, to enhance model capacity and precision, we scale the latent set size to 2048, reduce the learning rate to 5e-5, and halve the batch size to 960 for another 100k iterations, ensuring robust adaptation to higher-dimensional data spaces while maintaining computational efficiency.

## 4 Step1X-3D Texture Generation

Once the untextured 3D geometry is generated using the Step1X-3D framework, texture synthesis is performed via a multi-stage pipeline, as shown in Fig. 5. First, the raw geometry undergoes post-processing to ensure topological consistency and structural integrity (Sec. 4.1). Then, we prepare 3D assets for texture generation (Sec. 4.2). Next, a multi-view image generation model is fine-tuned on high-quality 3D datasets, incorporating geometric guidance via normal and position maps (Sec. 4.3). Finally, the generated multi-view images are super-resolved to 2048×2048 resolution before UV baking, followed by inpainting to complete the texture maps (Sec. 4.4).

### 4.1 Geometry Postprocess

To achieve high-fidelity texturing, we perform post-processing on the mesh geometry generated by the preceding geometric generation pipeline. The optimization process primarily employs the trimesh toolkit [16]. Specifically, we first verify the watertightness of the initial mesh, implementing hole-filling algorithms where non-manifold geometry is detected. Subsequently, we apply a remeshing operation that subdivides each triangular face into four sub-faces while enforcing Laplacian surface smoothing constraints. This remeshing procedure ensures uniform topological distribution and minimizes UV seam artifacts. Finally, we utilize the xAtlas [52] parameterization framework to generate optimized UV coordinates, which are then integrated into the final mesh representation. This systematic refinement pipeline guarantees geometric robustness for subsequent texture mapping.

### 4.2 Texture dataset preparation

Compared to geometry generation, the texture generation component does not require millions of training samples but instead places higher demands on texture quality and aesthetic metrics. From the 320K Objaverse dataset cleaned in Sec. 3.1, we further curated 30K 3D assets for multi-view generative model training. Specifically, we rendered each object with blender to produce six views (front, back, left, right, top, and bottom), along with corresponding albedo, normal map, and position map outputs with  $768 \times 768$  resolution.

### 4.3 Geometry-guided Multi-view Images Generation

**Single View to Multi-view Generation** Given a single view image and target multi-view camera poses (the conditions are defined as  $\mathcal{C}$ ), we aim to generate consistent multi-view images with a diffusion model  $D_{MV}$ . In details, we formulate the single-view to multi-view diffusion process as:

$$I_{1,2,\dots,N} = D_{MV}(z^{MV}, \mathcal{C}), \quad (5)$$

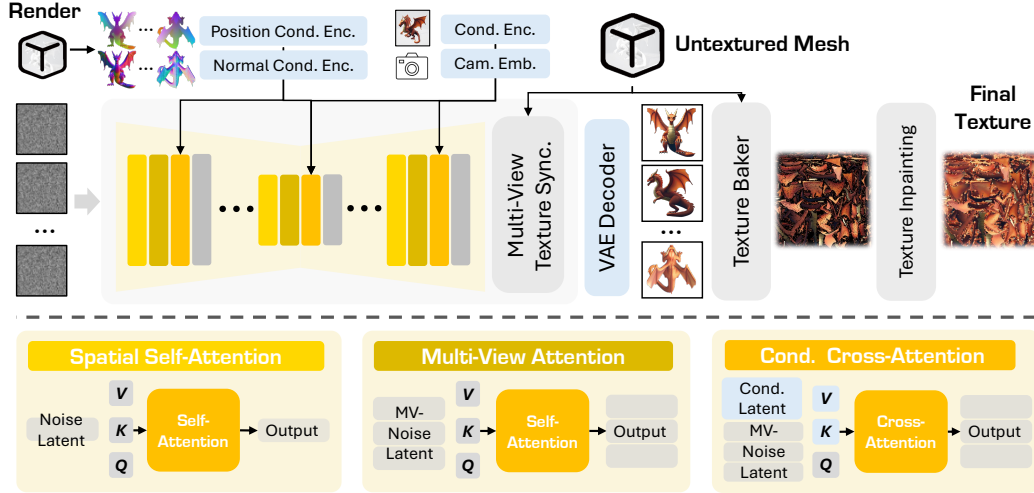


Figure 5: **Texturing module of Step1X-3D.**

where  $z^{MV}$  is a multi-view random noise. We use the pre-trained MV-Adapter [27] as the backbone to generate multi-view images with  $768 \times 768$  resolution and higher consistency. MV-Adapter exhibits two distinct advantages: the capability to generate high-resolution images and enhanced generalization performance. The high-resolution generation is primarily achieved through a memory-efficient epipolar-attention, which enables the production of  $768 \times 768$  resolution images under batch size constraints during the training process. The superior generalization capability stems from preserving the original spatial self-attention parameters of SD-XL [62], while introducing a triple parallelized attention architecture that simultaneously addresses generalization capacity, multi-view consistency, and conditional adherence. This design achieves an optimal balance between maintaining foundational model properties and acquiring specialized generation capabilities.

**Involve Geometry Guidance for Generation** Our objective is to generate reasonable and refined textures for the aforementioned untextured 3D meshes. To achieve this, during multi-view generation, in addition to conditioning on the provided single-view input, the injection of geometric information facilitates enhanced detail synthesis and improved alignment between textures and the underlying mesh surface. Specifically, we introduce two types of geometric guidance: normal maps and 3D position maps. The normal maps preserve fine-grained geometric details of the object, while the 3D position maps ensure accurate spatial correspondence between textures and mesh vertices across different viewpoints through 3D coordinate consistency. Both geometric representations are derived in the global world coordinate system. These geometric features are encoded via image-based encoders and subsequently injected into the backbone generation model through cross-attention mechanisms. This approach enables explicit geometric conditioning while maintaining the generative model’s ability to synthesize perceptually coherent textures.

**Synchronized Multi-view Images in Texture Domain** While the integration of cross-view attention and two geometric conditioning terms has achieved satisfactory multi-view consistency, inherent discrepancies between image space and UV space still introduce artifacts in synthesized textures, such as localized blurriness and discontinuous seams. To address this, we extend the MV-Adapter framework by introducing a texture-space synchronization module during inference. Unlike text-to-multi-view approaches like MVPaint [12] and SyncMVD [48] — which bypass explicit modeling of style references (sref) and content references (cref) between the input condition image and output multi-view images—our method eliminates the need for auxiliary refinement pipelines (e.g., Stable diffusion with controlNet) for multi-view synchronization. This design choice is justified by two considerations: 1) Our generator operates at a latent resolution of  $96 \times 96$ , which empirically provides sufficient texture representation capacity; 2) Joint optimization in a unified latent space inherently preserves texture coherence across views. Consequently, we implement texture synchronization exclusively through latent space alignment within a single diffusion model backbone, achieving parameter efficiency while maintaining visual fidelity.

In details, to predict the multi-view latent output  $z_{t+1}$ , we unproject the latent  $z_{i,t}$  of each view to texture space and  $\mathcal{P}$  represents UV mapping function. Then, we obtain the synchronized texture  $T'_{i,t}$  by fusing multi-view  $T_{i,t}$  and weight different views in texture space with the cosine similarities between the ray direction and per-pixel normal map. Further, the synchronized latent  $z'_{i,t}$  is obtained by projecting the  $T'_{i,t}$  with UV rasterization function  $\mathcal{P}$ . We formulate the whole procession of each denoise step as the following:

$$T_{i,t} = \mathcal{P}^{-1}(z_{i,t}), \quad (6)$$

$$T'_{i,t} = \sum_{i=1}^N \cos(v_i, n_i) * T_{i,t}, \quad (7)$$

$$z'_{i,t} = \mathcal{P}(T'_{i,t}), \quad (8)$$

$$z_{t+1} = D_{MV}(z'_t, \mathcal{C}). \quad (9)$$

#### 4.4 Bake Texture

Following conventional texture baking workflows [90, 12, 94], we adopt standard texture processing operations on multi-view projections of the object and reuse the texture baker tools in Hunyuan3D 2.0 [94]. First, multi-view images were upsampled to achieve a  $2048 \times 2048$  resolution and then are inversely projected onto the texture space. Due to occlusions and multi-view inconsistencies, this process inevitably introduces artifacts such as discontinuities and holes in the UV-mapped texture. To address this, we implement continuity-aware texture inpainting through iterative optimization, ensuring seamless texture synthesis across the entire surface. This post-processing stage effectively resolves topological ambiguities while preserving high-frequency texture details critical for photorealistic rendering.

## 5 Experiment

This section presents a comprehensive evaluation of the generation performance of Step1X-3D. First, we provide a detailed demonstration of Step1X-3D’s ability to generate geometry and texture conditioned on a single input image in Sec. 5.1. Next, we validate the model’s flexibility and controllability in Sec. 5.2. Finally, we conduct a thorough comparison between Step1X-3D and state-of-the-art (SOTA) methods, including both open-source (Trellis [85], Hunyuan3D 2.0 [71], and TripoSG [39]) and proprietary approaches (Meshy-4<sup>3</sup>, Tripo-v2.5<sup>4</sup>, and Rodin-v1.5<sup>5</sup>), across three key dimensions: quantitative metrics, user studies, and visual quality in Sec. 5.3.

### 5.1 The Visual Quality Results of Step1X-3D Assets

To evaluate Step1X-3D, we present the generated 3D assets from both geometric and textural dimensions in Fig. 6 and Fig. 7, respectively. To better showcase geometric details, we render multi-view normal maps from the generated meshes for 3D geometry visualization. As shown in Fig. 6, the first and sixth columns display input images, while the remaining columns present multi-view representations of different objects. Our test objects cover a wide variety of styles (cartoon, sketch, and photorealistic), geometric complexity (flat surfaces, hollow structures, and detail-rich objects), and spatial configurations (single objects and multi-object compositions). Across this diverse input images, Step1X-3D geometry generation model not only maintains strong similarity between the 3D mesh and input image, but also reconstructs plausible spatial structures for occluded regions with reasonable geometric details. These results demonstrate the crucial role of our specifically designed 3D diffusion model and VAE architecture in Step1X-3D, along with the significant improvement in generalization capability enabled by large-scale high-quality training data. Fig. 7 further demonstrates Step1X-3D’s texture generation capability through multi-view renders of textured 3D meshes. The texture generation model produces style-consistent textures across various input styles while maintaining high fidelity to the input image’s textural details. For occluded regions in the input image, by

<sup>3</sup>The API was invoked from Meshy platform in May 2025.

<sup>4</sup>The API was invoked from Tripo platform in May 2025.

<sup>5</sup>The API was invoked from Rodin platform in May 2025.

preserving the original SD-XL parameters and incorporating the target model’s normal maps and position maps as geometric guidance, Step1X-3D achieves plausible view completion with excellent multi-view consistency and precise geometry-texture alignment. In summary, Step1X-3D generates geometrically plausible 3D geometry with rich textures, where the final textured 3D meshes exhibit strong content and style matching with the input conditioning images.

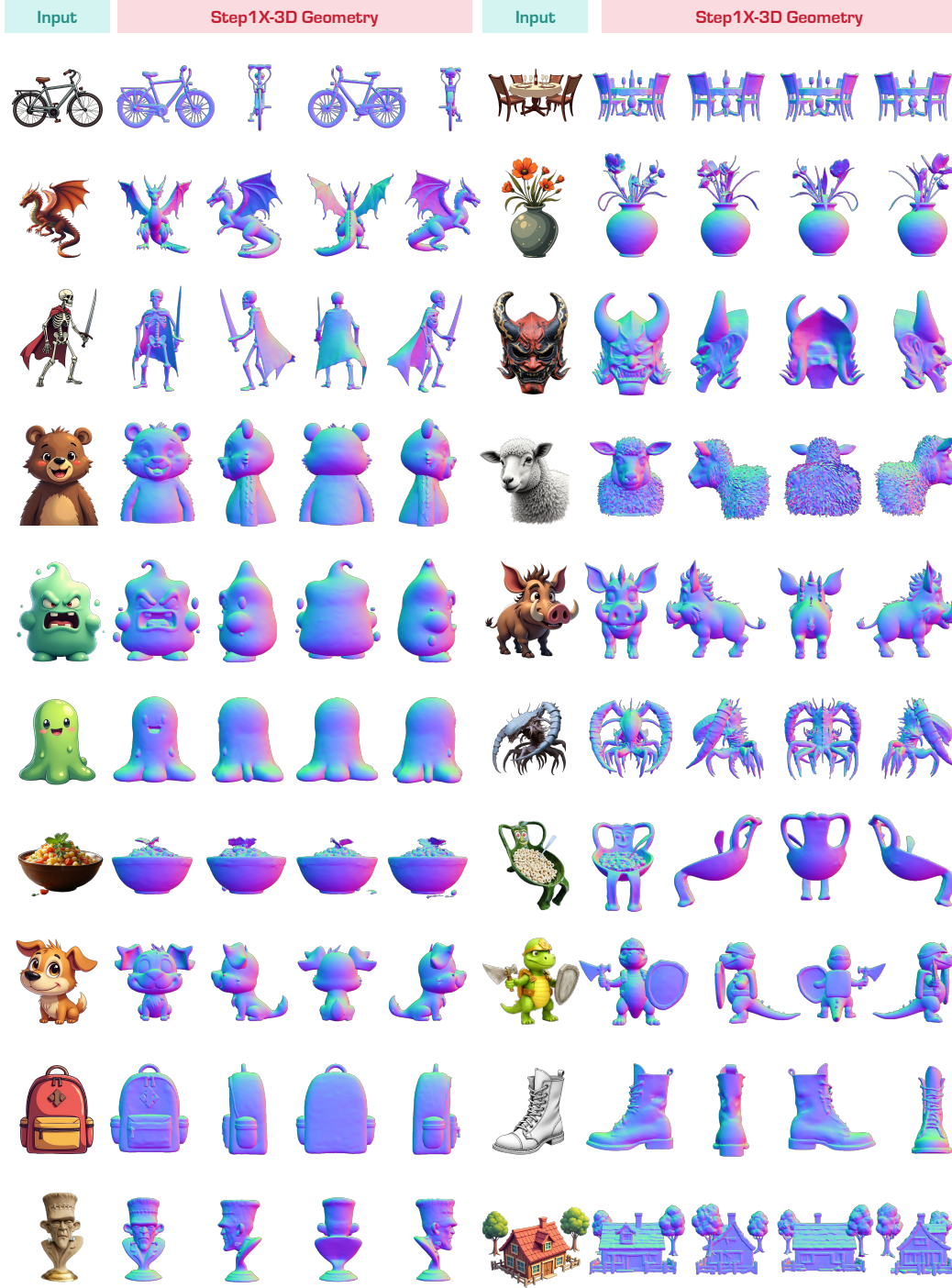


Figure 6: Samples of generated geometry by Step1X-3D.

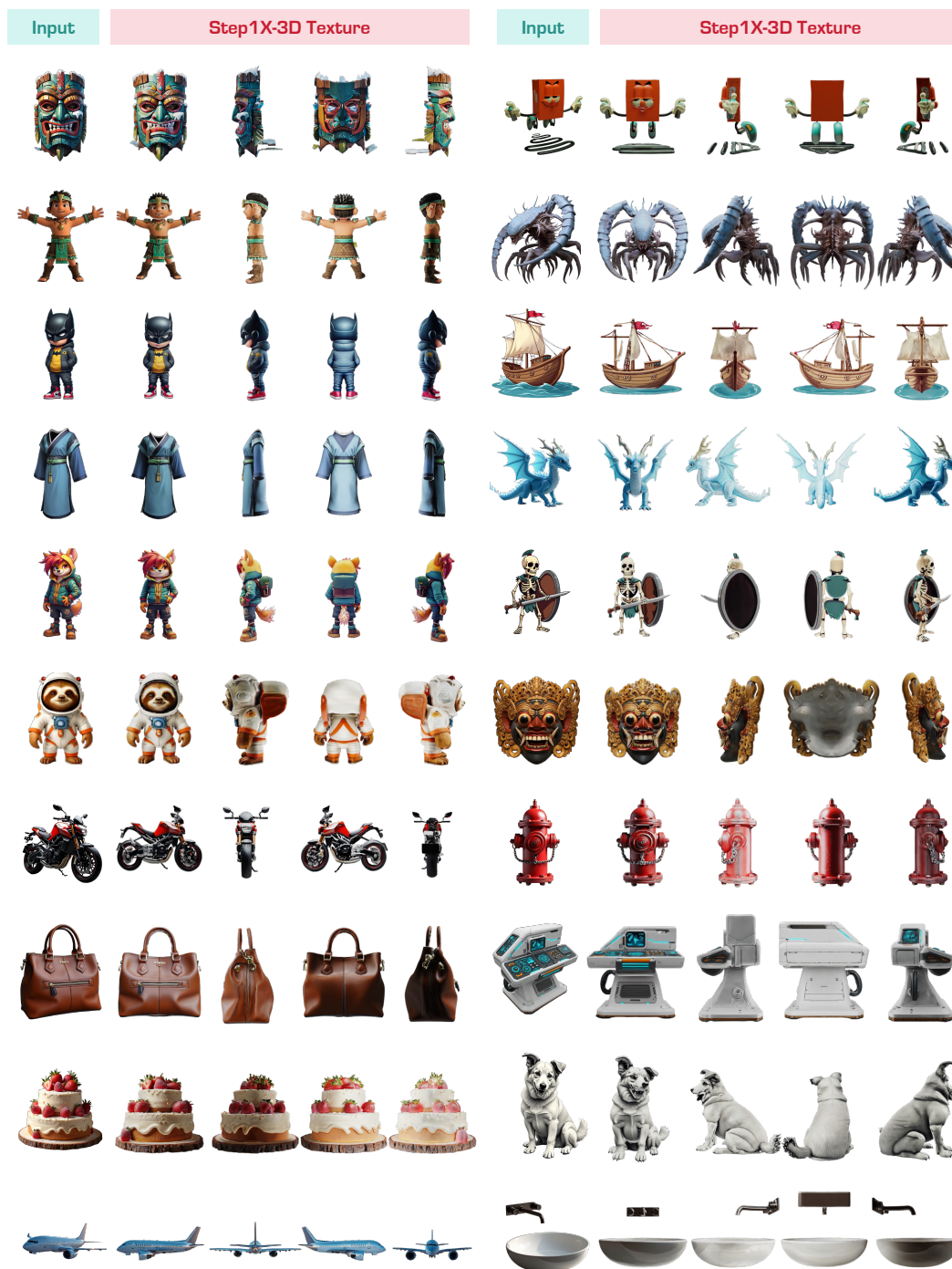


Figure 7: Samples of generated texture by Step1X-3D.

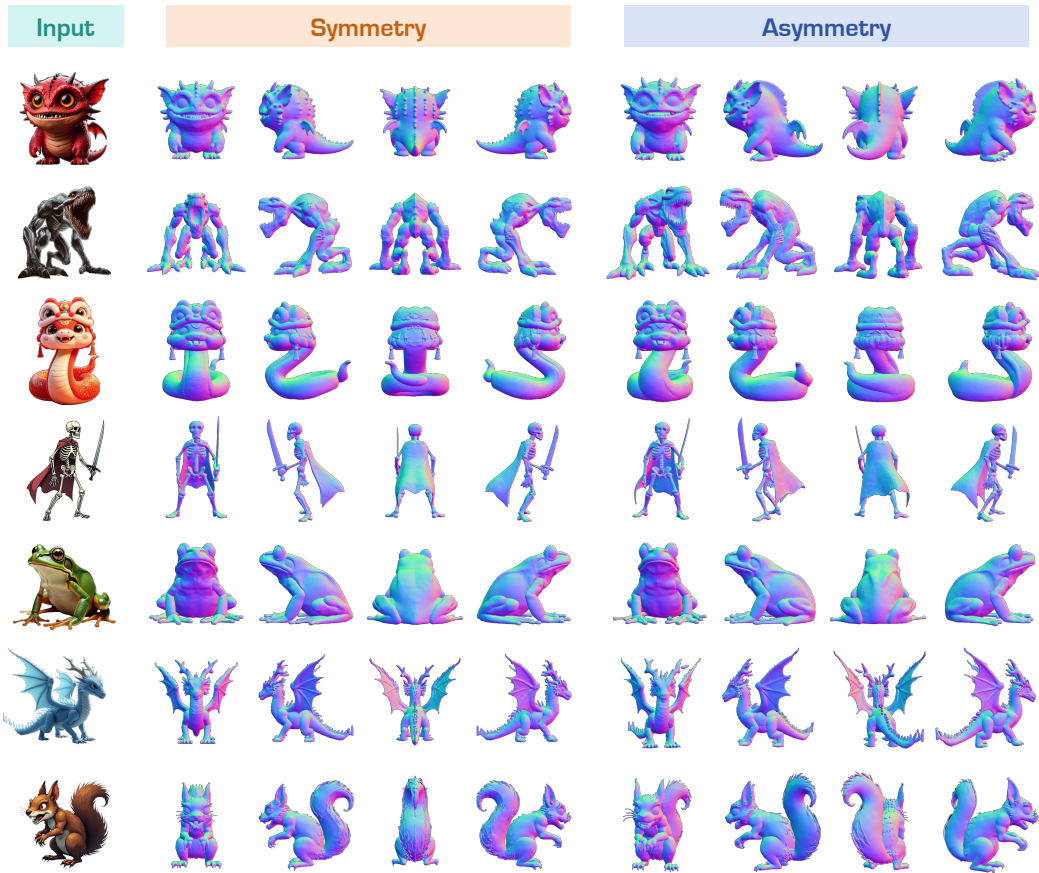


Figure 8: **Symmetry control.** Controllable 3D generation with symmetry or asymmetry geometry.

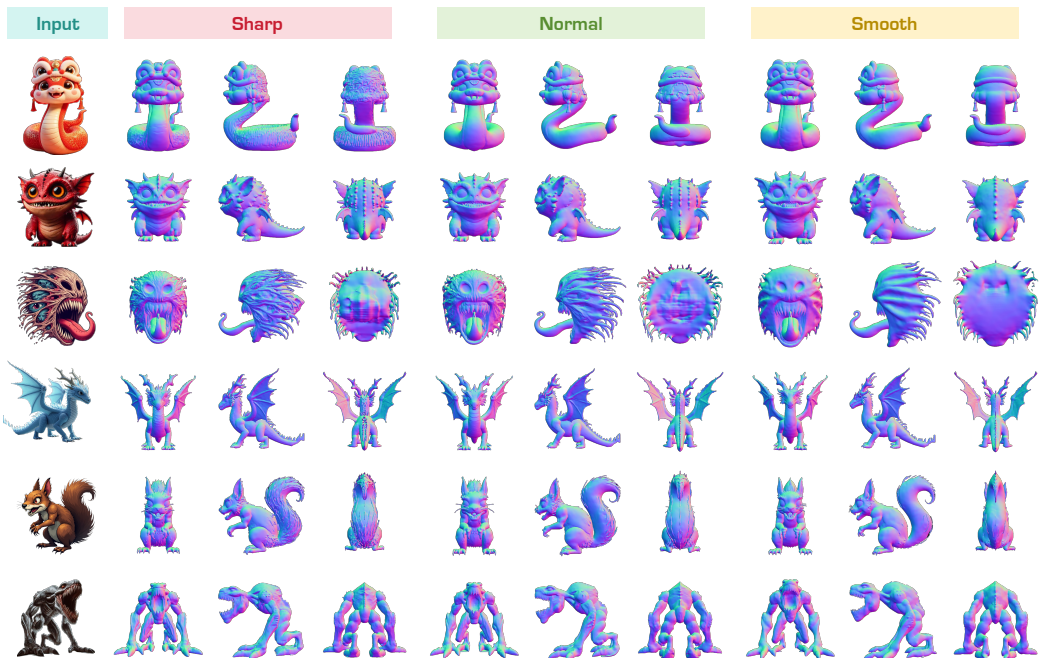


Figure 9: **Sharpness control.** Controllable 3D generation with sharp, normal or smooth geometry.



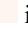
## 5.2 High Controllable 3D generation with LoRA Finetuning

Building upon a pre-trained geometry generation model for mesh reconstruction conditioned on a single image, we seamlessly integrate LoRA fine-tuning to enable flexible control over diverse 3D generation models. In this experiment, focusing on practical user applications, we specifically design two geometric control strategies: symmetry manipulation and hierarchical geometric detail adjustment. To implement these controls, we collected approximately 30,000 3D models and utilized the Step1O multimodal model to annotate each object based on (1) symmetry properties and (2) geometric detail levels (sharp, normal, and smooth). We show the high controllable 3D generation results in Fig. 8 and Fig. 9. To better capture geometric details in 3D meshes, we employ multi-view normal maps for geometric representation. Fig. 8 shows the results of geometry generation using "symmetry"/"asymmetry" captions. The first column shows the input image, columns 2–5 display four views (front, back, left, right) of the 3D object generated with a symmetry-conditioned caption, while columns 6–9 present the corresponding multi-view results for the asymmetry-conditioned generation. The results indicate that the generated 3D objects consistently adhere to their respective control instructions, with particularly pronounced compliance in the front and back views. Fig. 9 illustrates the hierarchical control of geometric details in detail. From left to right, we present the input conditioning image, followed by objects generated with the "sharp", "normal", and "smooth" labels respectively. Each object is represented using normal maps from front, right, and back views. Consistent with previous results, the generated objects demonstrate strong adherence to their corresponding control labels. This further validates the efficacy of Step1X-3D’s fine-tuning technique and demonstrates strong generalization capability in its geometry generation model.

## 5.3 Comparison Results with SOTA Methods

To further validate the effectiveness of Step1X-3D, we conducted comprehensive comparisons with existing state-of-the-art (SOTA) methods, including open-source approaches (Trellis [85], Hunyuan3D 2.0 [71], and TripoSG [39]) and proprietary systems (Tripo-v2.5, Rodin-v1.5, and Meshy-4). Specifically, we performed: (1) quantitative evaluations using both geometric and textural metrics; (2) user studies assessing perceived 3D quality through subjective scoring; (3) and visual comparisons of geometric and textural results across diverse input conditions.

**Quantity Comparisons and User Study across SOTA Methods** Beyond visual comparisons under diverse input conditions, we constructed a comprehensive benchmark dataset totaling 110 images in the wild. This benchmark incorporates: (1) example images from various 3D generation platforms (e.g., Tripo, Rodin and etc.), and (2) images generated by the Flux model covering 80 object categories from the COCO [40] dataset. Based on this test set, we systematically collected 3D assets generated by different methods for both quantitative evaluation and subjective user study.

Table 1: **Quantitative comparison on different models.** \* indicates closed-source models, and    indicate the best, second best, and third best performance respectively.

Model	Texture	Geometry		
	CLIP-Score $\uparrow$	Uni3D-I $\uparrow$	OpenShape <sub>sc</sub> -I $\uparrow$	OpenShape <sub>pb</sub> -I $\uparrow$
Rodin-v1.5*	0.845	0.270	0.104	0.081
Meshy-4*	0.796	0.357	0.137	0.099
Tripo-v2.5*	0.848	<b>0.366</b>	<b>0.140</b>	0.134
Trellis	0.848	0.353	0.136	<b>0.139</b>
Hunyuan3D 2.0	0.829	0.352	0.131	0.131
<b>Step1X-3D (Ours)</b>	<b>0.853</b>	0.361	0.139	0.130

We similarly designed quantitative metrics for both geometry and texture dimensions. For geometric evaluation, we leverage self-supervised multimodal models to perform feature matching between input 2D images and generated 3D point clouds (extracted from the output meshes). To ensure comprehensive and fair comparison, we utilized two distinct multimodal frameworks for feature extraction: Uni3D [96] and OpenShape [43], with cosine similarity serving as our similarity metric. For the OpenShape framework, which follows a self-supervised paradigm, we implemented both

SparseConv [13] and PointBERT [89] as backbone architectures. This yielded three distinct metrics for evaluating image-to-geometry alignment: Uni3D-I, OpenShape<sub>sc</sub>-I, and OpenShape<sub>pb</sub>-I, where higher scores indicate better geometric consistency with the input image. For texture evaluation, we adopted CLIP-Score [59] to measure semantic alignment. Specifically, we rendered multi-view images from textured 3D models at an elevation of 30° and azimuth angles of {0°, 90°, 180°, 270°} for semantic consistency assessment with input images. Quantitative results are presented in Tab. 1, with top and second-highest scores highlighted. Step1X-3D achieved the highest CLIP-Score and multiple second-highest rankings in geometric-semantic matching metrics. These superior results further demonstrate Step1X-3D’s robust generation capabilities.

We conducted a user study with 20 participants evaluating all 110 uncured test images. The assessment criteria for 3D models included: (1) geometric plausibility, (2) similarity to input images, (3) texture clarity, and (4) texture-geometry alignment. Participants rated each object on a 5-point Likert scale (1: lowest quality, 5: highest quality). As shown in Fig. 10, Step1X-3D achieves comparable performance to the current best-performing methods. However, we observe that all evaluated algorithms still fall significantly short of the theoretical upper bound, indicating considerable room for improvement before reaching production-ready quality. These findings underscore the importance of fostering open-source collaboration within the 3D research community to collectively advance the state-of-the-art.

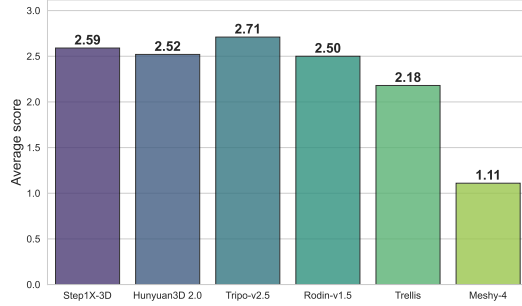


Figure 10: **User study.** Average user preference score is reported.

**Visual Comparisons across SOTA Methods** Fig. 11 and Fig. 12 present comparative results of geometric and textural outputs across different methods. Unlike previous visual comparisons, we address pose inconsistencies in the generated 3D meshes by implementing a unified evaluation protocol: (1) aligning both untextured and textured models in Unreal Engine for consistent pose normalization, and (2) compositing multiple objects into a single rendered image for direct comparison. This standardized approach reveals that Step1X-3D achieves comparable or superior performance relative to the best available methods.

## 6 Conclusion

Step1X-3D advances 3D generation by introducing an open-source, high-fidelity framework that decouples geometry and texture synthesis. Through rigorous data curation (2M assets) and a hybrid VAE-DiT architecture, it achieves promising results while enabling 2D-to-3D control transfer. We will release the models, training code, as well as the training data (excluding self-collected assets) to bridge the gap between proprietary and open research, fostering community progress toward production-ready 3D generation.

## 7 Limitations

Currently, we convert mesh to TSDF with grid resolution  $256^3$ . In future work, we will increase the grid resolution to achieve more accurate geometric details. Meanwhile, for the texture component, our current implementation is limited to albedo generation. We plan to extend this pipeline to support input image relighting and physically based rendering (PBR) material texture generation.

## 8 Contributors

**Algorithm contributors:** Weiyu Li, Xuanyang Zhang, Zheng Sun, Di Qi, Hao Li, Wei Cheng, Weiwei Cai, Zeming Li, Gang Yu, Xiangyu Zhang, Daxin Jiang

**Data contributors:** Shihao Wu, Jiarui Liu, Zihao Wang, Xiao Chen, Feipeng Tian, Jianxiong Pan

**Corresponding authors:** Xuanyang Zhang (zhangxuanyang@stepfun.com), Gang Yu (yugang@stepfun.com), Daxin Jiang (djiang@stepfun.com), Ping Tan (pingtan@ust.hk)



Figure 11: Qualitative comparison with SOTA methods on generated geometry.



Figure 12: Qualitative comparison with SOTA methods on generated texture.

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Raphael Bensadoun, Tom Monnier, Yanir Kleiman, Filippos Kokkinos, Yawar Siddiqui, Mahendra Kariya, Omri Harosh, Roman Shapovalov, Benjamin Graham, Emilien Garreau, et al. Meta 3d gen. *arXiv preprint arXiv:2407.02599*, 2024.
- [3] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- [4] Mario Botsch, Leif Kobbelt, Mark Pauly, Pierre Alliez, and Bruno Lévy. *Polygon mesh processing*. CRC press, 2010.
- [5] 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.
- [6] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015.
- [7] Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner. Text2tex: Text-driven texture synthesis via diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 18558–18568, 2023.
- [8] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 22246–22256, 2023.
- [9] Rui Chen, Jianfeng Zhang, Yixun Liang, Guan Luo, Weiyu Li, Jiarui Liu, Xiu Li, Xiaoxiao Long, Jiashi Feng, and Ping Tan. Dora: Sampling and benchmarking for 3d shape variational auto-encoders. *arXiv preprint arXiv:2412.17808*, 2024.
- [10] Sijin Chen, Xin Chen, Anqi Pang, Xianfang Zeng, Wei Cheng, Yijun Fu, Fukun Yin, Billzb Wang, Jingyi Yu, Gang Yu, et al. Meshxl: Neural coordinate field for generative 3d foundation models. *Advances in Neural Information Processing Systems*, 37:97141–97166, 2024.
- [11] Yiwen Chen, Tong He, Di Huang, Weicai Ye, Sijin Chen, Jiaxiang Tang, Xin Chen, Zhongang Cai, Lei Yang, Gang Yu, et al. Meshanything: Artist-created mesh generation with autoregressive transformers. *arXiv preprint arXiv:2406.10163*, 2024.
- [12] Wei Cheng, Juncheng Mu, Xianfang Zeng, Xin Chen, Anqi Pang, Chi Zhang, Zhibin Wang, Bin Fu, Gang Yu, Ziwei Liu, et al. Mvpaint: Synchronized multi-view diffusion for painting anything 3d. *arXiv preprint arXiv:2411.02336*, 2024.
- [13] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3075–3084, 2019.
- [14] Jasmine Collins, Shubham Goel, Kenan Deng, Achleshwar Luthra, Leon Xu, Erhan Gundogdu, Xi Zhang, Tomas F Yago Vicente, Thomas Dideriksen, Himanshu Arora, Matthieu Guillaumin, and Jitendra Malik. Abo: Dataset and benchmarks for real-world 3d object understanding. *CVPR*, 2022.
- [15] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers, 2024.
- [16] Dawson-Haggerty et al. trimesh.
- [17] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. Objaverse-xl: A universe of 10m+ 3d objects. *arXiv preprint arXiv:2307.05663*, 2023.
- [18] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13142–13153, 2023.
- [19] Ding Ding, Zeqian Ju, Yichong Leng, Songxiang Liu, Tong Liu, Zeyu Shang, Kai Shen, Wei Song, Xu Tan, Heyi Tang, et al. Kimi-audio technical report. *arXiv preprint arXiv:2504.18425*, 2025.

- [20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [21] Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binqiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture. *International Journal of Computer Vision*, pages 1–25, 2021.
- [22] Google Gemini2. Experiment with gemini 2.0 flash native image generation, 2025.
- [23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [24] 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.
- [25] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- [26] Ailin Huang, Boyong Wu, Bruce Wang, Chao Yan, Chen Hu, Chengli Feng, Fei Tian, Feiyu Shen, Jingbei Li, Mingrui Chen, et al. Step-audio: Unified understanding and generation in intelligent speech interaction. *arXiv preprint arXiv:2502.11946*, 2025.
- [27] Zehuan Huang, Yuan-Chen Guo, Haoran Wang, Ran Yi, Lizhuang Ma, Yan-Pei Cao, and Lu Sheng. Mv-adapter: Multi-view consistent image generation made easy. *arXiv preprint arXiv:2412.03632*, 2024.
- [28] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [29] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In *International conference on machine learning*, pages 4651–4664. PMLR, 2021.
- [30] Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 867–876, 2022.
- [31] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*, 2023.
- [32] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023.
- [33] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- [34] Zeqiang Lai, Yunfei Zhao, Zibo Zhao, Haolin Liu, Fuyun Wang, Huiwen Shi, Xianghui Yang, Qinxiong Lin, Jinwei Huang, Yuhong Liu, Jie Jiang, Chunchao Guo, and Xiangyu Yue. Unleashing vecset diffusion model for fast shape generation, 2025.
- [35] Houyi Li, Wenzhen Zheng, Jingcheng Hu, Qiufeng Wang, Hanshan Zhang, Zili Wang, Shijie Xuyang, Yuantao Fan, Shuigeng Zhou, Xiangyu Zhang, et al. Predictable scale: Part i—optimal hyperparameter scaling law in large language model pretraining. *arXiv preprint arXiv:2503.04715*, 2025.
- [36] Peng Li, Yuan Liu, Xiaoxiao Long, Feihu Zhang, Cheng Lin, Mengfei Li, Xingqun Qi, Shanghang Zhang, Wei Xue, Wenhan Luo, et al. Era3d: high-resolution multiview diffusion using efficient row-wise attention. *Advances in Neural Information Processing Systems*, 37:55975–56000, 2024.
- [37] Wei Yu Li, Rui Chen, Xuelin Chen, and Ping Tan. Sweetdreamer: Aligning geometric priors in 2d diffusion for consistent text-to-3d. *arxiv:2310.02596*, 2023.
- [38] Wei Yu Li, Jiarui Liu, Rui Chen, Yixun Liang, Xuelin Chen, Ping Tan, and Xiaoxiao Long. Craftsman: High-fidelity mesh generation with 3d native generation and interactive geometry refiner. *arXiv preprint arXiv:2405.14979*, 2024.
- [39] Yangguang Li, Zi-Xin Zou, Zexiang Liu, Dehu Wang, Yuan Liang, Zhipeng Yu, Xingchao Liu, Yuan-Chen Guo, Ding Liang, Wanli Ouyang, et al. Triposg: High-fidelity 3d shape synthesis using large-scale rectified flow models. *arXiv preprint arXiv:2502.06608*, 2025.
- [40] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision—ECCV 2014: 13th European conference, zurich, Switzerland, September 6–12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014.
- [41] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

- [42] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [43] Minghua Liu, Ruoxi Shi, Kaiming Kuang, Yin hao Zhu, Xuanlin Li, Shizhong Han, Hong Cai, Fatih Porikli, and Hao Su. Openshape: Scaling up 3d shape representation towards open-world understanding. *Advances in neural information processing systems*, 36:44860–44879, 2023.
- [44] Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming Wang, Honghao Fu, Chunrui Han, et al. Step1x-edit: A practical framework for general image editing. *arXiv preprint arXiv:2504.17761*, 2025.
- [45] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- [46] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology, limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402.17177*, 2024.
- [47] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Sync-dreamer: Generating multiview-consistent images from a single-view image. In *The Twelfth International Conference on Learning Representations*, 2024.
- [48] Yuxin Liu, Minshan Xie, Hanyuan Liu, and Tien-Tsin Wong. Text-guided texturing by synchronized multi-view diffusion. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024.
- [49] 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. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9970–9980, 2024.
- [50] 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.
- [51] Guoqing Ma, Haoyang Huang, Kun Yan, Liangyu Chen, Nan Duan, Shengming Yin, Changyi Wan, Ranchen Ming, Xiaoni Song, Xing Chen, et al. Step-video-t2v technical report: The practice, challenges, and future of video foundation model. *arXiv preprint arXiv:2502.10248*, 2025.
- [52] Markus Worchel et al. xatlas.
- [53] A.L.F. Meister. *Generalia de genesi figurarum planarum et inde pendentibus earum affectionibus*. 1769.
- [54] 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. *Communications of the ACM*, 65(1):99–106, 2021.
- [55] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022.
- [56] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- [57] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023.
- [58] 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*, 2023.
- [59] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. Pmlr, 2021.
- [60] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021.
- [61] Elad Richardson, Gal Metzger, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-guided texturing of 3d shapes. In *ACM SIGGRAPH 2023 conference proceedings*, pages 1–11, 2023.
- [62] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [63] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.

- [64] Aditya Sanghi, Hang Chu, Joseph G Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshan. Clip-forge: Towards zero-shot text-to-shape generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18603–18613, 2022.
- [65] 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. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19615–19625, 2024.
- [66] 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. In *European Conference on Computer Vision*, pages 1–18. Springer, 2024.
- [67] Jiaxiang Tang, Zhaoshuo Li, Zekun Hao, Xian Liu, Gang Zeng, Ming-Yu Liu, and Qinsheng Zhang. Edgerunner: Auto-regressive auto-encoder for artistic mesh generation. *arXiv preprint arXiv:2409.18114*, 2024.
- [68] Shitao Tang, Jiacheng Chen, Dilin Wang, Chengzhou Tang, Fuyang Zhang, Yuchen Fan, Vikas Chandra, Yasutaka Furukawa, and Rakesh Ranjan. Mvdifffusion++: A dense high-resolution multi-view diffusion model for single or sparse-view 3d object reconstruction. In *European Conference on Computer Vision*, pages 175–191. Springer, 2024.
- [69] Zhenyu Tang, Junwu Zhang, Xinhua Cheng, Wangbo Yu, Chaoran Feng, Yatian Pang, Bin Lin, and Li Yuan. Cycle3d: High-quality and consistent image-to-3d generation via generation-reconstruction cycle. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 7320–7328, 2025.
- [70] Tencent Hunyuan3D Team. Hunyuan3d 1.0: A unified framework for text-to-3d and image-to-3d generation, 2024.
- [71] Tencent Hunyuan3D Team. Hunyuan3d 2.0: Scaling diffusion models for high resolution textured 3d assets generation, 2025.
- [72] 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.
- [73] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [74] Vikram Voleti, Chun-Han Yao, Mark Boss, Adam Letts, David Pankratz, Dmitry Tochilkin, Christian Laforte, Robin Rombach, and Varun Jampani. Sv3d: Novel multi-view synthesis and 3d generation from a single image using latent video diffusion. In *European Conference on Computer Vision*, pages 439–457. Springer, 2024.
- [75] Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wenten Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.
- [76] Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al. Neural codec language models are zero-shot text to speech synthesizers. *arXiv preprint arXiv:2301.02111*, 2023.
- [77] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J. Guibas. Normalized object coordinate space for category-level 6d object pose and size estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [78] Zhengyi Wang, Jonathan Lorraine, Yikai Wang, Hang Su, Jun Zhu, Sanja Fidler, and Xiaohui Zeng. Llama-mesh: Unifying 3d mesh generation with language models. *arXiv preprint arXiv:2411.09595*, 2024.
- [79] 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. In *European Conference on Computer Vision*, pages 57–74. Springer, 2024.
- [80] Si-Tong Wei, Rui-Huan Wang, Chuan-Zhi Zhou, Baoquan Chen, and Peng-Shuai Wang. Octgpt: Octree-based multiscale autoregressive models for 3d shape generation. *arXiv preprint arXiv:2504.09975*, 2025.
- [81] Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan, Valentin Deschaintre, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Meshlm: Large reconstruction model for high-quality meshes. *arXiv preprint arXiv:2404.12385*, 2024.

- [82] Haohan Weng, Zibo Zhao, Biwen Lei, Xianghui Yang, Jian Liu, Zeqiang Lai, Zhuo Chen, Yuhong Liu, Jie Jiang, Chunchao Guo, et al. Scaling mesh generation via compressive tokenization. *arXiv preprint arXiv:2411.07025*, 2024.
- [83] Kailu Wu, Fangfu Liu, Zhihan Cai, Runjie Yan, Hanyang Wang, Yating Hu, Yueqi Duan, and Kaisheng Ma. Unique3d: High-quality and efficient 3d mesh generation from a single image. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [84] Shuang Wu, Youtian Lin, Feihu Zhang, Yifei Zeng, Jingxi Xu, Philip Torr, Xun Cao, and Yao Yao. Direct3d: Scalable image-to-3d generation via 3d latent diffusion transformer. *arXiv preprint arXiv:2405.14832*, 2024.
- [85] Jianfeng Xiang, Zelong Lv, Sicheng Xu, Yu Deng, Ruicheng Wang, Bowen Zhang, Dong Chen, Xin Tong, and Jiaolong Yang. Structured 3d latents for scalable and versatile 3d generation. *arXiv preprint arXiv:2412.01506*, 2024.
- [86] 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.
- [87] Yinghao Xu, Zifan Shi, Wang Yifan, Hansheng Chen, Ceyuan Yang, Sida Peng, Yujun Shen, and Gordon Wetzstein. Grm: Large gaussian reconstruction model for efficient 3d reconstruction and generation. In *European Conference on Computer Vision*, pages 1–20. Springer, 2024.
- [88] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.
- [89] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert: Pre-training 3d point cloud transformers with masked point modeling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19313–19322, 2022.
- [90] Xianfang Zeng, Xin Chen, Zhongqi Qi, Wen Liu, Zibo Zhao, Zhibin Wang, Bin Fu, Yong Liu, and Gang Yu. Paint3d: Paint anything 3d with lighting-less texture diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4252–4262, 2024.
- [91] Biao Zhang, Jiapeng Tang, Matthias Niessner, and Peter Wonka. 3dshape2vecset: A 3d shape representation for neural fields and generative diffusion models. *ACM Transactions On Graphics (TOG)*, 42(4):1–16, 2023.
- [92] Longwen Zhang, Ziyu Wang, Qixuan Zhang, Qiwei Qiu, Anqi Pang, Haoran Jiang, Wei Yang, Lan Xu, and Jingyi Yu. Clay: A controllable large-scale generative model for creating high-quality 3d assets. *ACM Transactions on Graphics (TOG)*, 43(4):1–20, 2024.
- [93] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023.
- [94] Zibo Zhao, Zeqiang Lai, Qingxiang Lin, Yunfei Zhao, Haolin Liu, Shuhui Yang, Yifei Feng, Mingxin Yang, Sheng Zhang, Xianghui Yang, et al. Hunyuan3d 2.0: Scaling diffusion models for high resolution textured 3d assets generation. *arXiv preprint arXiv:2501.12202*, 2025.
- [95] 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*, 36:73969–73982, 2023.
- [96] Junsheng Zhou, Jinsheng Wang, Baorui Ma, Yu-Shen Liu, Tiejun Huang, and Xinlong Wang. Uni3d: Exploring unified 3d representation at scale. *arXiv preprint arXiv:2310.06773*, 2023.
- [97] Cailin Zhuang, Yaoqi Hu, Xuanyang Zhang, Wei Cheng, Jiacheng Bao, Shengqi Liu, Yiyang Yang, Xianfang Zeng, Gang Yu, and Ming Li. Styleme3d: Stylization with disentangled priors by multiple encoders on 3d gaussians. *arXiv preprint arXiv:2504.15281*, 2025.
- [98] Zi-Xin Zou, Zhipeng Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Yan-Pei Cao, and Song-Hai Zhang. Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10324–10335, June 2024.
- [99] Qi Zuo, Xiaodong Gu, Lingteng Qiu, Yuan Dong, Weihao Yuan, Rui Peng, Siyu Zhu, Liefeng Bo, Zilong Dong, Qixing Huang, et al. Videomv: Consistent multi-view generation based on large video generative model. 2024.