

Progressive Text-to-3D Generation for Automatic 3D Prototyping

HAN YI, Department of Computer Science, University of North Carolina at Chapel Hill, North Carolina, USA

ZHEDONG ZHENG, FST and ICI, University of Macau, Macau, China

XIANGYU XU, Xi'an Jiaotong University, Shaanxi, China

TAT-SENG CHUA, National University of Singapore, Singapore, Singapore

The challenge of text-to-3D generation lies in accurately and efficiently crafting 3D objects based on natural language descriptions, a capability that promises substantial reduction in manual design efforts and offers an intuitive interface for user interaction with digital environments. Despite recent advancements, effective recovery of fine-grained details and efficient optimization of high-resolution 3D outputs remain critical hurdles. Drawing inspiration from the efficacious paradigm of progressive learning, we present a novel Multi-scale Triplane Network (MTN) architecture coupled with a tailored progressive learning strategy. As the name implies, the Multi-scale Triplane Network consists of four triplanes transitioning from low to high resolution. This hierarchical structure allows the low-resolution triplane to serve as an initial shape for the high-resolution counterparts, easing the inherent complexity of the optimization process. Furthermore, we introduce the progressive learning scheme that systematically guides the network to shift its attention from prominent coarse-grained structures to intricate fine-grained patterns. This strategic progression ensures that the focus of the model evolves towards emulating the subtlest aspects of the described 3D object. Our experiment verifies that the proposed method performs favorably against contemporary methods. Even for the complex and nuanced textual descriptions, our method consistently excels, delivering robust and viable 3D shapes where other methods falter.

ACM Reference Format:

Han Yi, Zhedong Zheng, Xiangyu Xu, and Tat-Seng Chua. 2026. Progressive Text-to-3D Generation for Automatic 3D Prototyping. *ACM Trans. Multimedia Comput. Commun. Appl.* 9, 4 (March 2026), 19 pages. <https://doi.org/10.1145/nnnnnnnnnnnnnnnn>

1 INTRODUCTION

Designing digital models for manufacturing [11, 13] is often time-consuming and labor-intensive. To streamline this process, researchers are developing more intuitive methods for 3D object generation, such as using text prompts (see Figure 1). The aim of the text-to-3D generation task is to automatically create a 3D object draft from a natural description, thus cutting down the design efforts from the ground up.

In recent years, text-to-3D generation has reported rapid development due to the breakthrough of text-to-image diffusion models [9, 41, 48]. For instance, the pioneer work DreamFusion [43] leverages the 2D Stable Diffusion and proposes Score Distillation Sampling (SDS) algorithm to generate a variety of 3D objects using only text prompts. However, there remain two problems:

Authors' addresses: Han Yi, Department of Computer Science, University of North Carolina at Chapel Hill, North Carolina, USA; Zhedong Zheng, FST and ICI, University of Macau, 999078, Macau, China; Xiangyu Xu, Xi'an Jiaotong University, 310013, Shaanxi, China; Tat-Seng Chua, National University of Singapore, Singapore, Singapore.

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

© 2009 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 1551-6857/2026/3-ART

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

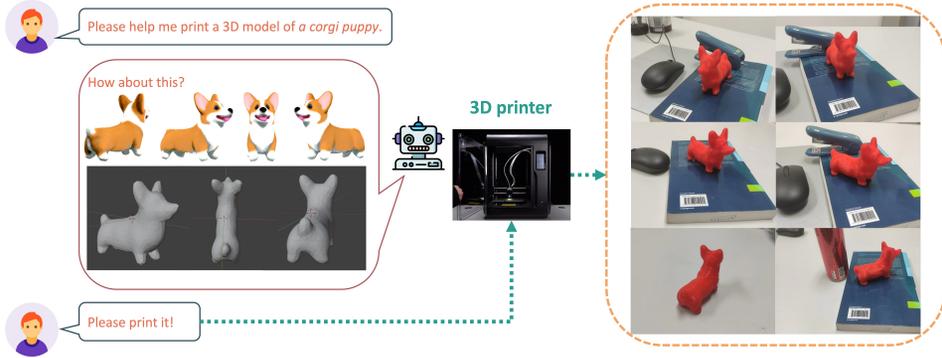


Fig. 1. The proposed algorithm facilitates effortless and interactive creation of high-quality 3D meshes from natural language descriptions, which can then be utilized for 3D printing. The six images at the right show the corresponding physical 3D printed model from multiple perspectives. Our output meshes are ready for 3D printing. (We add the book and mouse as the size reference.)

1) The inherent optimization complexity of 3D high-resolution objects. It is hard to directly map one sentence to one high-dimension 3D object, especially in the form of Neural Radiance Fields (NeRF) [37]. This leads to either generation collapse or extended training duration for model convergence. 2) Lack of fine-grained details. We notice that some works report blurred results [35, 43, 55]. This is due to the use of a fixed training strategy, *i.e.*, focusing on global fidelity all the time while ignoring local parts.

In an attempt to overcome the above-mentioned challenges, we propose a progressive text-to-3D generation model that can gradually refine details to produce high-quality 3D objects (see Figure 2). 1) For the first problem, we introduce a novel network structure, namely, Multi-scale Triplane Network (MTN) consisting of four triplanes ranging from low to high resolution. In the initial phases of training, we sample low-resolution features from the corresponding low-resolution triplane to capture the basic global geometric shape. As training advances, we fix the former low-resolution triplanes and gradually shift our focus to triplanes with a higher resolution. Such a progressive structure facilitates the model to capture different-level features in a step-by-step manner and thus enhances the geometric and textural nuances of the 3D model, such as color and texture. 2) For the second problem, we adopt a progressive learning strategy focusing on two key factors, *i.e.*, time step t and camera radius. In particular, unlike existing 2D diffusion models that utilize random sampling, we adopt a large t during the initial stages to guide the global structure. As the training progresses, we transition to a smaller t to refine visual details. Meanwhile, we gradually adjust the radius of the camera to approach the object more closely. This enables the camera to initially focus on capturing the global structure and later shift its attention to the local details. Our contributions are as follows:

- We introduce a Multi-scale Triplane Network (MTN) to effectively tackle the challenge of text-to-3D generation in a bottom-up manner. This hierarchical structure progresses from rough to fine-grained details, leveraging initial low-resolution shapes to streamline the high-resolution optimization, overcoming complexities faced by prior methods.
- We propose a progressive learning strategy tailored for the Multi-scale Triplane Network. It simultaneously reduces the camera radius and time step t in diffusion to refine details of the 3D model in a coarse-to-fine manner, ensuring superior capture of subtle details in the generated 3D models.

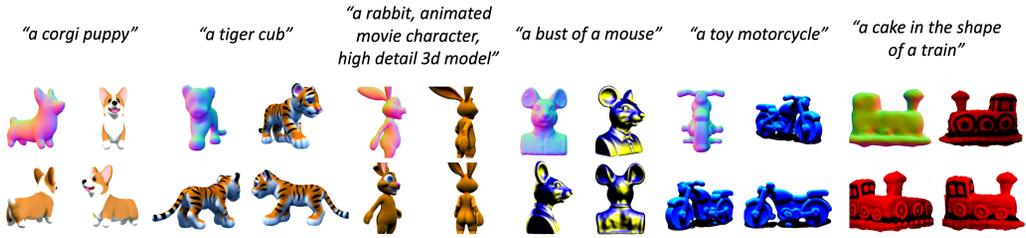


Fig. 2. Our method is able to generate high-quality 3D outputs from various text prompts using the proposed Multi-scale Triplane Network (MTN). We display both mesh normals and the generated results obtained from texts of varying lengths. Specifically, our approach showcases the ability to create animal meshes and industrial products. Moreover, automatic color rendering is applied when a common color is applicable for such a category.

- Albeit simple, extensive experiments show that the proposed method could achieve high-resolution outputs that align closely with natural language descriptions. We expect this work to pave the way for automatic 3D printing via intuitive human-machine interaction.

2 RELATED WORK

3D Generative Modeling The realm of 3D generative modeling has seen extensive exploration across diverse representation types, including voxel grids [25, 29, 52], point clouds [32, 53, 68], meshes [14–17, 33, 40, 45], implicit fields [8, 59, 60, 65], triplanes [7], and octrees [19]. While many traditional approaches hinge on 3D assets as training data, the challenge of acquiring such data at scale has spurred alternative strategies. Addressing the inherent challenge of obtaining 3D assets for training, some recent endeavors have turned to 2D supervision. Leveraging ubiquitous 2D images, models, *e.g.*, pi-GAN [4], EG3D [3], MagicMirror [67], and GIRAFFE [42] have supervised 2D renderings of 3D models through adversarial loss against 2D image datasets. While these approaches hold potential, a recurring challenge is that they are often restricted to specific domains, *e.g.*, human faces [21, 62] and bodies [66], limiting their versatility and hindering expansive creative freedom in 3D design. In our study, we pivot towards text-to-3D generation, with the goal of crafting visually favorable 3D objects guided by diverse text prompts.

Text-to-3D Generation The success of text-to-image generation models has driven substantial progress in the emerging field of text-to-3D object generation. Notably, the integration of CLIP into models, *e.g.*, CLIP-forge [46], Dream Fields [20], Text2Mesh [36], CLIPmesh [38], and CLIP-NeRF [54] has been a significant advancement. These approaches harness CLIP to optimize 3D representations, ensuring that 2D renderings resonate with textual prompts. A defining advantage of such techniques is their ability to bypass the need for costly 3D training data, though a trade-off in terms of the realism of the resultant 3D models has been observed. More recent advancements, such as DreamFusion [43], which proposes Score Distillation Sampling (SDS) Loss, SJC [55], Magic3D [28], and Latent-NeRF [35], have showcased the merits of employing robust text-to-image diffusion models as a robust 2D prior, elevating the quality and realism of text-to-3D generation. Such a visual prior, capitalizing on the potential of diffusion models, has led to outcomes with higher fidelity and diversity, as well as reduced generation time. Along this line, Fantasia3D [5] employs disentangled modeling of geometry and appearance, enhancing fidelity and realism while offering better control over both properties. Meanwhile, ProlificDreamer [57] introduces Variational Score Distillation (VSD) Loss, serving as a replacement for SDS Loss. This enhancement has resulted in outputs characterized by higher resolution and increased diversity in 3D representations. Despite these advances, the multi-face (Janus) problem remains. To address this, Zero-1-to-3 [31],

Image-Dream [56], and MVDream [47] introduce multi-view diffusion models trained on extensive multi-view datasets to ensure multi-view consistency. Additionally, Bidiff [10] presents a unified framework integrating 3D and 2D diffusion processes to preserve both 3D fidelity and 2D texture richness. While substantial, these contributions differ from our focus on enhancing 3D representation quality and can complement our method. Triplane-based methods, such as Instant3D [24], DIRECT-3D [30], and TPA3D [58] represent a promising alternative within the NeRF-based text-to-3D landscape. By leveraging efficient Triplane representations, these approaches achieve a balance between computational efficiency and output quality. These methods reveal the potential of Triplane representations to elevate text-to-3D generation tasks and align closely with the principles of our approach. Recently, 3D Gaussian Splatting [22] has emerged as an alternative to NeRF. Methods like DreamGaussian [51], GSGEN [6], GaussianDreamer [63], and LucidDreamer [27] have applied this representation to text-to-3D generation. Though faster, these approaches often compromise the high quality characteristic of NeRF-based methods and require post-processing to convert Gaussian representations into NeRF or meshes, adding computational overhead. Therefore, we focus on NeRF-based methods for their superior quality and fidelity. Furthermore, while large-scale diffusion-based pipelines such as Step1X-3D [26] and Hunyuan3D-2.0 [64] have further pushed the fidelity of text-to-3D generation by leveraging >1B-parameter backbones, these methods rely on heavy pre-training and industrial-scale compute, operating under a very different assumption than single-GPU training regimes. Building upon the principles of high-fidelity, NeRF-based generation, our method specifically targets the triplane representation. Unlike prior triplane-based approaches that either use triplanes as a static reconstruction target (e.g., Instant3D [24]) or integrate them within a single-stage generator (e.g., TPA3D [58]), our method introduces a multi-resolution triplane/trivector field jointly optimized within a diffusion-guided coarse-to-fine learning pipeline, enabling progressive refinement of geometry and texture in a unified framework.

3 METHOD

3.1 Multi-scale Triplane

An overview of our Multi-scale Triplane Network (MTN) is shown in Figure 3 (a). In particular, MTN is composed of four triplanes [3] ranging from low to high resolutions. Each triplane leverages three axis-aligned 2D feature planes $\mathbf{F}_{xy}^m, \mathbf{F}_{xz}^m, \mathbf{F}_{yz}^m \in \mathbb{R}^{N_m \times N_m \times C}$, $m = 1, 2, 3$. N_m denotes spatial resolution, while C is the dimension of the channels and m represents the training stage. Note that a large N_m results in a substantial GPU memory cost. Therefore, for the last triplane, we essentially employ a trivector instead to optimize memory usage and support higher resolution. This trivector configuration leverages the axis-aligned vectors $\mathbf{F}_x^4, \mathbf{F}_y^4, \mathbf{F}_z^4 \in \mathbb{R}^{N_4 \times 1 \times C}$ with a resolution of $N_4 \times 1$ and C .

Given any 3D coordinate point $p \in \mathbb{R}^3$, we project this coordinate onto each of these orthogonal feature planes and sample feature vectors via interpolation. We then sum these three vectors $f^m(p) = \mathbf{F}_{xy}^m(p) + \mathbf{F}_{xz}^m(p) + \mathbf{F}_{yz}^m(p)$ for $m = 1, 2, 3$ as position features for the first three triplanes, while $f^4(p) = \mathbf{F}_x^4(p) + \mathbf{F}_y^4(p) + \mathbf{F}_z^4(p)$ for the last trivector. To aggregate multi-scale features, we further fuse the different level position features together as $h^m(p) = \sum_{k=1}^m (f^k(p))$. After obtaining the multi-scale representation, we follow [49] to transform the summed position features into the Fourier domain. Subsequently, the Fourier features are fed forward into a lightweight triplane decoder to estimate color and density [37]. We deploy a Multi-Layer Perceptron (MLP) as the triplane decoder. Finally, to calculate the loss, we apply neural volume rendering techniques [37] to project the 3D representation onto an RGB image I , which is the input of the Diffusion model.

Discussion. Why is a multi-scale structure crucial? As shown in Figure 3, we apply triplanes with different resolutions to capture features at multiple scales. This approach is designed to mimic

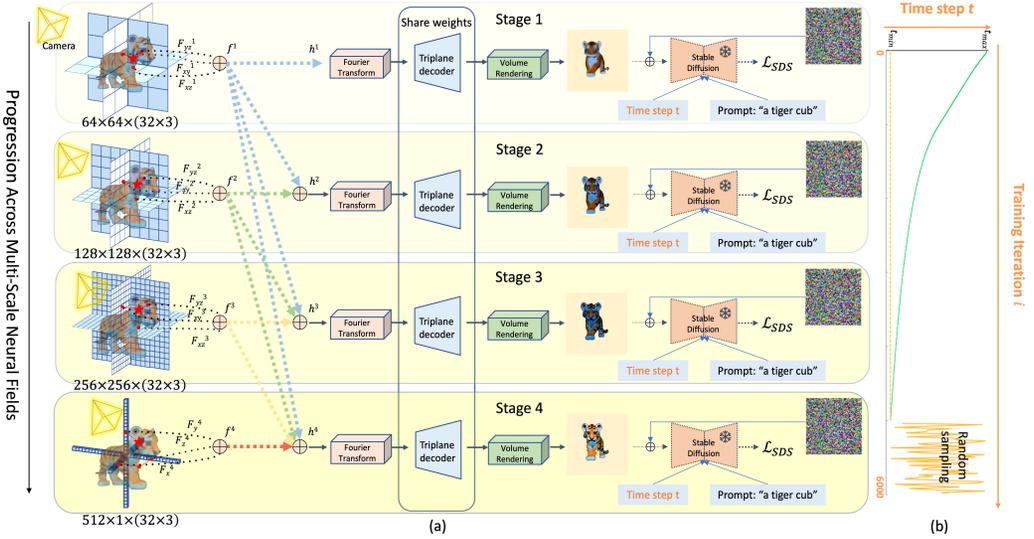


Fig. 3. Overview of the proposed Multi-scale Triplane Network (MTN). **(a)** Given the text prompts, e.g., “a tiger cub”, MTN generates 3D representations using Multi-scale Neural Fields, utilizing four triplanes varying in resolution. To save memory costs and enable the highest resolution, we make a trade-off to deploy the high-dimension trivector format as the triplane alternative. First, by casting rays from a random camera position and view, we can sample a lot of 3D points along each ray and then encode their corresponding features by projecting them onto triplanes. After the 3D input encoding, the network uses a Fourier transform, a triplane decoder, and volume rendering. The Fourier feature transform [49] enables the triplane decoder to learn high-frequency information. The network employs Fourier transform, a shallow MLP triplane decoder, and volume rendering to convert the 3D representation into RGB images. Training progresses in four stages, starting with low-resolution triplanes for global geometric insights, and gradually shifting to higher-resolution triplanes for detailed refinement. **(b)** Concurrently, as training proceeds, the time step t undergoes progressive adjustments, and the camera also approaches the neural field progressively, emphasizing the refinement of local features. To update the parameters, we employ a frozen Stable Diffusion model to estimate the injected noise on the rendered image (e.g., tiger) and then backpropagate the gradient.

the human recognition system, which transitions from recognizing basic elements to more intricate details when observing 3D objects. For example, when a person sees a new object, they first perceive its overarching structure and then refine the details through foveal vision. During the early stages of training, we extract low-resolution features from the corresponding low-resolution triplane. Each point on the low-resolution triplane, obtained through interpolation from a coarse grid, encompasses a broader field of view, providing global geometric insights. As training progresses, we gradually shift our focus to higher-resolution triplanes, which can capture intricate features and refine details such as subtle shading and texture nuances. This process facilitates the optimization of high-scale features, especially when low-scale features have already been well-optimized. This multi-scale approach is conceptually similar to curriculum learning [2], where the model starts with simpler tasks and gradually advances to more complex ones. In the experiments, we observe that the proposed method achieves visual enhancements in both shape and texture of the model, even for complex descriptions.

Optimization objective. Given the projected image I , we apply Score Distillation Sampling (SDS) [43] to distill 2D image priors from the pretrained 2D diffusion model ϵ_ϕ . The loss on 2D projection is then back-propagated to update differentiable 3D representations. In particular, the proposed 3D model can be typically depicted as a parametric function $I = g_\theta(P)$, where I represents

the images produced at distinct camera poses, and P is the set of multiple positions p . Here, g denotes the volumetric rendering mechanism, and θ embodies a coordinate-based MLP and triplanes that portray a 3D scene. To estimate the projection quality, we adopt the pretrained diffusion model, which is well aligned with text prompts y . The one-time denoising forward can be formulated as $\epsilon_\phi(I_t; y, t)$ to predict the noise ϵ given the noisy image I_t , time step t , and text embedding y . Therefore, the gradient of the SDS loss can be formulated as:

$$\nabla_{\theta} \mathcal{L}_{SDS}(\phi, g_{\theta}(P)) = \mathbb{E}_{t, \epsilon} \left[(\epsilon_{\phi}(I_t; y, t) - \epsilon) \frac{\partial I_t}{\partial \theta} \right],$$

where ϵ is a noise term following a standard normal distribution and I_t denotes the noisy image. Following the setting in the diffusion model [9, 41, 48], the noisy image can be formulated as a linear process $I_t = \sqrt{\bar{\alpha}_t}I + \sqrt{1 - \bar{\alpha}_t}\epsilon$, where $\bar{\alpha}_t$ is a predefined time-dependent constant. Besides, it is worth noting that the diffusion model parameter ϕ is frozen. The purpose of this denoising function is to offer the text-aware guidance to update θ . If the projection I is well-aligned with the text y , the noise on I_t is easy to predict. Otherwise, we will punish the 3D model.

3.2 Progressive Learning Strategy

Another essential element underlying the proposed method is the employment of a progressive learning strategy, focusing on two critical parameters, *i.e.*, the time step t and camera radius.

Progressive time step sampling. We first introduce a progressive time step (t) sampling approach. It is motivated by the observation that the default uniform t -sampling in SDS training often results in inefficiencies and inaccuracies due to the broad-range random sampling. Our approach, therefore, emphasizes a gradual reduction of the time step, directing the model to transition from coarse to detailed learning (See Figure 3 (b)). In the early phases of training, we adopt larger time steps to add a substantial amount of noise into the image. During the noise recovery process, the network is driven to focus on the low-frequency global structure signal. As training advances and the global structure stabilizes, we decrease to smaller time steps. In this stage, the network is demanded to recover the high-frequency fine-grained pattern according to the context. It facilitates the model in refining local details, such as textures and shades. We define the rate of change of variable t as:

$$\frac{dt}{di} = \beta v(t), \quad (1)$$

where $v(t)$ controls how t changes with respect to the training iteration i and is manually designed. β is a positive constant. We define $v(t)$ piece-wise:

$$v(t) = \begin{cases} -\exp(\frac{t-n_2}{m_2}) & \text{if } t > n_2 \\ -1.0 & \text{if } n_1 \leq t \leq n_2 \\ -\exp(\frac{t-n_1}{m_1}) & \text{if } t < n_1, \end{cases} \quad (2)$$

Here, $v(t) < 0$ implies $\frac{dt}{di} < 0$, indicating that t decreases as training progresses. Our design ensures that t decreases rapidly at the beginning ($t > n_2$), linearly in the middle ($n_1 \leq t \leq n_2$), and more mildly towards the end ($t < n_1$). After the time step t decreases to t_{\min} , we revert to random sampling from a uniform distribution as: $t \sim \mathcal{U}(t_{\min}, t_{\max})$, where $\mathcal{U}(t_{\min}, t_{\max})$ denotes uniform sampling within the interval from t_{\min} to t_{\max} . It reintroduces randomness to maintain the vibrancy of the coloration of the 3D model. We notice that a concurrent work, Dreamtime [18], also employs a similar non-increasing t -sampling strategy. However, such a strategy sometimes tends to overfit the local details, and inadvertently change the global illumination. Therefore, it is crucial to avoid the consistent use of extremely small time steps at the end of training. Different

from Dreamtime [18], our method decreases t with the training step at a much steeper pace and employs a mixture of both deterministic and random sampling as shown in Figure 3 (b).

Progressive radius. Simultaneously, our approach also incorporates a dynamic camera radius considering the camera movements in the real world. Typically, eyes will move closer for detailed object observation. Motivated by this behavior, we dynamically adjust the camera radius during the multi-scale learning. During the low-scale triplane stage, which focuses on broader geometric structures, we utilize a large camera radius to cover the entire object. As we move to the high-scale triplane stage, which refines local model details, the camera radius is reduced to closely focus on finer details of the 3D scene. This progressive radius strategy is intuitive and directly impacts resolution, aiding in feature learning across varying scales. In the ablation study, we also verify the effectiveness of this strategy (See Section 4.3).

3.3 Implementation Details

Neural field rendering structure. The proposed MTN consists of three triplanes and one trivector varying in resolution. The resolutions of the triplanes $N_1, N_2, N_3 = 64, 128, 256$, and the number of channels $C = 32$. For the trivector, we set $N_4 = 512$. During the Neural Field optimization, camera positions are randomly sampled in spherical coordinates. The azimuth angles, polar angles and fovy range are randomly sampled between $[-180^\circ, 180^\circ]$, $[45^\circ, 105^\circ]$, and $[10^\circ, 30^\circ]$, respectively. For spherical radius of the camera, the initial $R \in [3.0, 3.5]$ and gradually decreases to $R \in [1.8, 2.1]$. **In practice, we use the full radius range during the first 3,000 iterations, reduce it to 0.8 \times between 3,000–4,000 iterations, 0.7 \times between 4,000–5,000 iterations, and 0.6 \times thereafter.**

Prompts. For prompt augmentation, the default view-dependent prompt augmentation appends corresponding view, e.g., “front view”, “back view”, and “side view” according to the camera position. However, we adopt the strategy from Perp-Neg [1], leveraging geometric properties to enhance the diffusion model’s alignment with user prompts. This approach enriches original prompts with view-dependent conditional text embeddings based on sampled camera positions, ensuring the rendered image adheres to the desired view. Specifically, if the azimuth angle $\phi \in [-90^\circ, 90^\circ]$, a soft embedding is interpolated between “front view” and “side view” based on ϕ and appended to the original text embedding. Conversely, for $\phi \notin [-90^\circ, 90^\circ]$, the algorithm interpolates between “back view” and “side view” embeddings. This nuanced addition ensures more accurate and user-aligned renderings.

Diffusion model. We deploy DeepFloyd-IF [23] as the guidance model to provide 2D image priors. For time step (t) sampling in SDS, the Stable-DreamFusion uses random sampling $t \sim \mathcal{U}(20, 980)$. In our proposed approach, the time step t is set to decrease from 980 to 20. Through a grid search, we empirically set an optimal prior weight configuration as $\{m_1 = 50, m_2 = 150, n_1 = 500, n_2 = 800\}$ to control the rate of decrease. Following existing works [1, 28, 43], we also adopt the viewpoint-aware prompts by appending prompts such as “front view”, “side view”, and “back view”.

Optimization. The number of total iterations is 6000 and the batch size is 1. We employ the Adan optimizer [61] with learning rate of 1×10^{-3} , weight decay of 2×10^{-5} . Following the existing work [3], we apply two regularization terms, i.e., TV regularization and L2 regularization, to prevent floating clouds. The model can converge within one hour on a V100 GPU. Specifically, we configure the training process with 3,000 iterations for the first stage, followed by 1,000 iterations each for the second, third, and final stage, respectively.

4 EXPERIMENT

In this section, we assess the capability of our method to produce high-fidelity 3D objects according to natural language prompts. We primarily consider three key evaluation aspects: (1) alignment with the text, particularly focusing on key words in the sentence; (2) intricate texture

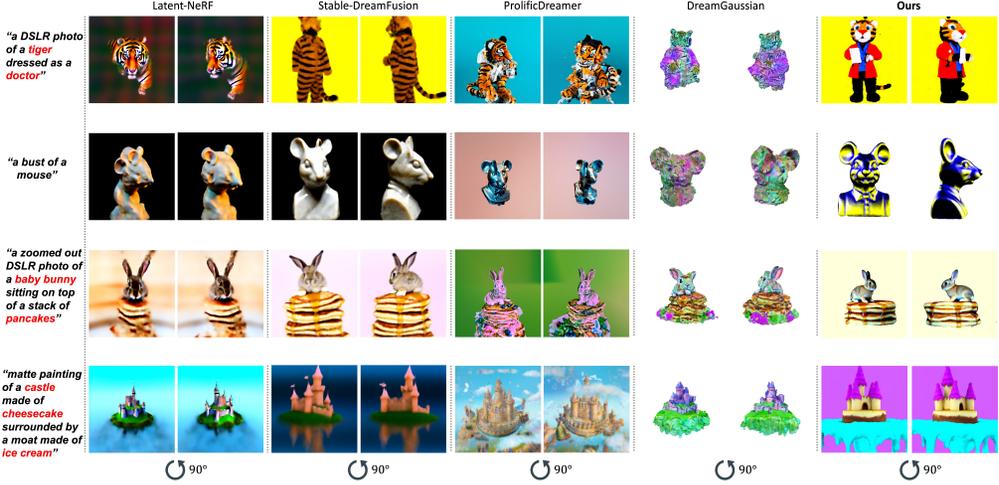


Fig. 4. Qualitative comparisons for text-to-3D generation among our method, Latent-NeRF [35], Stable-DreamFusion [50], ProlificDreamer [57], and DreamGaussian [51]. Here we show the 2D projection of the front view and side view of the 3D model. We observe that the proposed method could generate a higher-fidelity 3D representation aligned with the given description, reducing the extra post-processing costs. **In the last row, despite ProlificDreamer [57] and Latent-NeRF [35] achieves good visual quality, they generally miss the keyword “cheesecake” and “ice cream”.**

details; and (3) consistent geometric shape, especially in localized parts, *e.g.*, ears and tails. Due to the space limitation, we mainly compare our approach against four widely-used text-to-3D frameworks. Since DreamFusion [43] is not publicly available, we utilize the open-source variant, Stable-DreamFusion [50]. Besides, we also compare the proposed method with other three competitive works, *i.e.*, Latent-NeRF [35], ProlificDreamer [57], and DreamGaussian [51].

4.1 Qualitative Evaluation

As shown in Figure 4, we could observe that our method outperforms prior competitive approaches in terms of text alignment, texture details, and geometric precision. The qualitative analysis reveals the superior performance of our method in generating realistic and accurate 3D representations aligned with textual prompts. In the first row, we observe notable deficiencies in Latent-NeRF [35], which struggles to produce a coherent 3D model. While Stable-DreamFusion [50] manages to generate a tiger avatar, it fails to incorporate the crucial keyword “doctor”. ProlificDreamer [57], despite its high output resolution, erroneously includes unrelated elements, such as a camera, on the tiger’s hand, which is obviously inconsistent with the specified theme of “a tiger doctor.” DreamGaussian [51], on the other hand, successfully identifies the “tiger face” element but falters in rendering the rest of the model, resulting in an overall geometry that appears unconventional. In contrast, our proposed method seamlessly integrates the textual cues to craft a detailed representation of a tiger doctor, complete with a book in its hands. In the second row, our method presents a refined geometric shape with correct shading on the bust, surpassing Stable-DreamFusion [50], which erroneously places a tail on the head. Similarly, the outputs from Latent-NeRF [35], ProlificDreamer [57], and DreamGaussian [51] display inaccuracies in head shape, notably featuring three ears and multi-face. Additionally, DreamGaussian [51] shows discrepancies in color saturation, resulting in outputs that are excessively vibrant. Simultaneously, our method distinguishes itself by depicting nuanced features such as the necktie and buttons on the mouse. In the third row, our method accurately

captures the keyword “baby bunny”, showcasing a natural geometric shape with clear edges and appropriate features. Conversely, both Latent-NeRF [35] and Stable-DreamFusion [50] continue to struggle with the multi-face and multi-ear issue. ProlificDreamer [57], and DreamGaussian [51], while offering high-resolution outputs, fall short in aligning their geometric shapes and color fidelity with the textual prompt, underscoring the critical balance between resolution and semantic coherence. In the last row, our method aligns well with the given text prompt, accurately capturing the three keywords “castle”, “cheesecake”, and “ice cream”, and generates high-quality 3D outputs with exquisite textures. In contrast, other methods primarily focus on the keyword “castle” and overlook the additional critical details. Although ProlificDreamer [57] produces a visually appealing scene with diverse features, its output appears foggy and cloud-filled, which deviates noticeably from the given prompt. To further validate the robustness of our method, Figures 5 present additional qualitative comparisons, including recently proposed methods such as GSGEN [6] and LucidDreamer [27], evaluated across a wide range of prompts. Both GSGEN and LucidDreamer suffer from the multi-head (Janus) artifact, leading to inconsistencies in multi-view rendering—an issue our method effectively overcomes. In summary, our method consistently excels in generating reliable, semantically aligned, and visually coherent 3D models. The outputs not only adhere closely to the textual prompts but also exhibit intuitive geometric structures that align well with human perception.

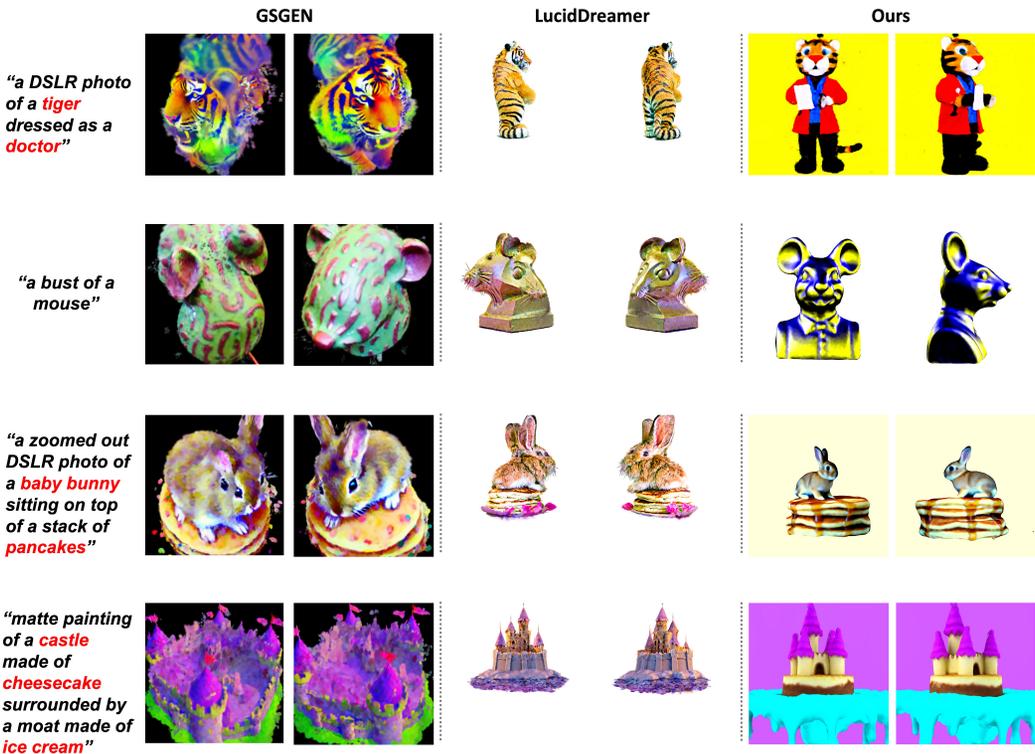


Fig. 5. Qualitative comparisons with GSGEN and LucidDreamer.

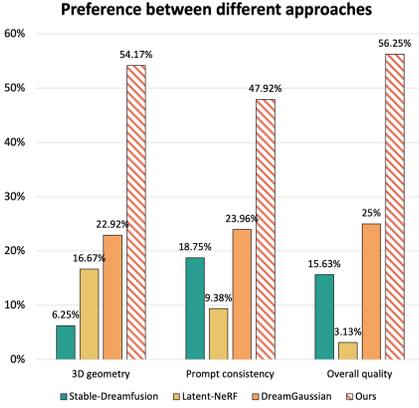


Fig. 6. User study on visual quality. The proposed method excels in 3D geometry, closely aligns with user prompts, and outperforms two competitive approaches in overall quality.

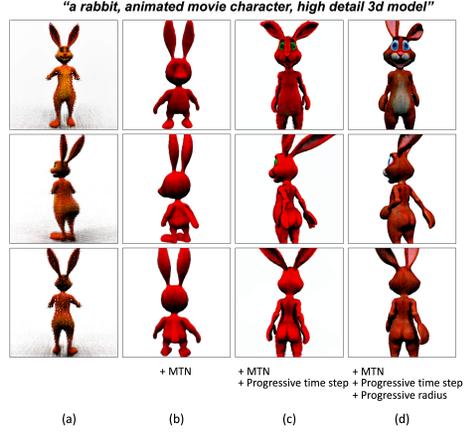


Fig. 7. Ablation study of the primary components. (a) Single triplane; (b) add MTN architecture; (c) add Progressive time step strategy; (d) add Progressive radius, which is our full method. Our full model crafts a delicate geometric shape and achieves accurate texture.

User Study. For a more comprehensive evaluation, we conduct a user study with 96 participants. We evaluate our model against three prevailing and basic approaches, *e.g.*, Latent-NeRF [35], Stable-DreamFusion [50], and DreamGaussian [51] in three key aspects: 3D geometry, prompt consistency, and overall quality. We randomly select 96 prompts from the standard set of 153 prompts and generate 3D models, using Stable-DreamFusion [50], Latent-NeRF [35], DreamGaussian [51], and our approach. Participants are then asked to rank the models based on the aforementioned criteria. As shown in Figure 6, our visual results outperform other methods across multiple metrics, attracting preferences from 56.25% of participants for overall quality, 54.17% for 3D geometry, and 47.92% for prompt consistency. This highlights the efficacy of our approach across various evaluation criteria.

4.2 Quantitative Evaluation

Since our task is a generation problem, we lack 3D ground-truth meshes for direct quantitative comparison of differences. Therefore, we follow the existing work, *i.e.*, DreamFusion [43], to evaluate the alignment between 2D projected images and the text prompt. In particular, we adopt the CLIP R-Precision [44] to evaluate the retrieval performance for both RGB images and depth maps. The RGB images serve as an indicator of texture quality, while the depth maps represent the geometric shape. A higher score indicates better performance. This evaluation is conducted using three pre-trained CLIP models with different model sizes, *i.e.*, CLIP B/32, CLIP B/16, and CLIP L/14. For a fair comparison, we also adopt 153 standard prompts from Dreamfields [20]. As shown in Table 1, we observe that our

Table 1. Quantitative comparisons with competitive methods. The best precision in every column is in **bold**. We do not include ProlificDreamer [57] in this table, since it is extremely time-consuming, requiring about 11 hours per prompt for just the first training stage.

Method	R-Precision (%) \uparrow					
	CLIP B/32		CLIP B/16		CLIP L/14	
	RGB	DEPTH	RGB	DEPTH	RGB	DEPTH
GT images	77.1	-	79.1	-	-	-
Latent-NeRF	48.4	37.1	52.9	40.6	59.5	40.9
Stable-Dreamfusion	56.4	45.9	60.3	45.8	58.3	42.9
DreamGaussian	61.3	48.7	61.9	49.2	61.7	45.8
Ours	62.6	53.1	62.6	51.9	64.8	47.6

Table 2. The ablation study investigates the impact of different components, with the best precision highlighted in **bold** for each column. The ablation study validates the effectiveness of the proposed MTN architecture, progressive time step, and progressive radius. Notably, the full model (MTN-full) achieves the highest level of text-visual semantic alignment.

Method	MTN	Progressive Time Step	Progressive Radius	R-Precision (%) \uparrow							
				CLIP B/32		CLIP B/16		CLIP L/14		Mean	
				RGB	DEPTH	RGB	DEPTH	RGB	DEPTH	RGB	DEPTH
Single triplane				46.8	38.4	51.8	41.1	53.9	41.4	50.8	40.3
MTN	✓			57.8	46.7	58.2	46.2	62.2	42.8	59.4	45.2
MTN-t	✓	✓		60.2	52.7	61.2	51.0	63.5	43.5	61.6	49.1
MTN-r	✓		✓	57.9	48.5	60.4	48.8	62.4	42.7	60.2	46.7
MTN-full	✓	✓	✓	62.6	53.1	62.6	51.9	64.8	47.6	63.3	50.9

method consistently achieves the highest R-Precision scores across all three metrics in terms of both RGB texture and depth, indicating a significant advantage.

4.3 Ablation Study and Further Discussion

Effectiveness of Multi-scale Triplanes. We first investigate the impact of the multi-scale triplane architecture to substantiate its advantages. As shown in Table 2, we could observe that the multi-scale architecture facilitates both texture and geometric shape learning. Specifically, the RGB R-Precision is improved with a large margin +8.6% on average, while the mean depth R-Precision increases +4.9%. We also provide a visualization result in Figure 7 (b). The basic single-scale triplane structure results in a 3D output that misses intricate details both texturally and geometrically, evident in incomplete hands, tails, and the presence of floating points. (Noted that for the single-triplane baseline, we use a resolution of 512×512 , the same as the final resolution in our progressive multi-scale approach. This ensures that the single-triplane setup has comparable capacity to represent high-resolution details, allowing for a meaningful comparison.) In contrast, the multi-scale network gradually leverages the multi-scale information, yielding a more smooth geometric shape with clear edges. While there are still imperfections, the rabbit now possesses a more complete form, particularly noticeable in its overall silhouette.

Further analysis is presented in Figure 8 (a), which shows renderings with different combinations of triplanes. For example, the first image uses only the 64-triplane, while the second adds the 128-triplane, and so on. For the final triplane, we use a trivector to optimize memory usage and support higher resolutions. The total number of training iterations is the same for all four combinations in Figure 8 (a). This illustrates how the generation quality improves as higher-scale triplanes are added. Sampling from higher-resolution triplanes enhances details and sharp edges. More detailed explanations are provided in the "Why is a multi-scale structure crucial?" section in our paper.

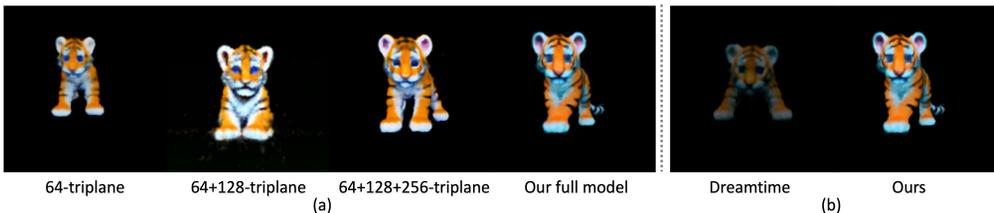


Fig. 8. (a) Hierarchical outputs. (b) Ablation on strategies.

Effectiveness of Progressive Learning. Here we further evaluate the impact of progressive time step sampling and progressive radius. (1) As shown in the third row of Table 2, the MTN with



Fig. 9. Visual results at different training iterations (3,000–6,000). The model progressively refines coarse geometric shapes into detailed and realistic textures, demonstrating the effectiveness of the coarse-to-fine learning strategy.

only progressive time step strategy could further improve the text alignment by +2.2% texture and +3.9% geometry quality on average. This is because the small time step towards the end of learning shifts the focus to high-frequency details, significantly improving the overall visual quality. (2) Similar to how humans often take a closer look to examine object details, our model, when applying the progressive radius approach, performs even better, showing a +1.7% improvement on the local texture details. As the camera gets closer, the 2D projection and the optimization objects both emphasize local quality, resulting in a refined 3D model. As a result, the culmination of these strategies leads to a final output that is both detailed and visually appealing (see Figure 7 (d)). Additionally, we visualize the training progression from 3000 to 6000 iterations (see Figure 9). Even at 3000 iterations, the model already reconstructs a coarse yet recognizable shape, while later iterations refine appearance and texture quality. This demonstrates the coarse-to-fine nature of our progressive strategy, which first captures global structure and then enhances local details.

Effectiveness of Random Sampling in the final stage. Unlike the concurrent DreamTime strategy [18], which does not use random sampling in the final stage, we adopt random sampling. As shown in Figure 8 (b), while keeping other hyperparameters consistent, our approach results in better illumination conditions and clearer edges compared to DreamTime [18].

Compatibility and Scalability. The proposed method is compatible with various pre-trained diffusion models as supervision, and can be easily extended to further improve the quality of generation. For instance, our approach can integrate seamlessly with the state-of-the-art multi-view diffusion model MVDream [47], which effectively tackles the multi-face problem by emphasizing multi-view consistency. The combination enables a superior 3D consistency and exquisite textures and verifies the compatibility and scalability of our method (see Figure 10). To further validate our method’s effectiveness, we compare the combination of MTN + MVDream with the original MVDream in Table 3. Our MTN outperforms the NeRF component used in MVDream in terms of R-Precision, while also requiring less training time and fewer parameters. It is important to note that we did not tune the hyperparameters for MTN, instead directly using those optimized for NeRF in the original MVDream, highlighting the robustness and plug-and-play nature of our method.

Table 3. Comparison of NeRF backbone.

NeRF Backbone	Diffusion	R-Precision (%) ↑	Training Time ↓	#Params ↓
NeRF (from MVDream)	MVDream	67.1	1.5 hours	12.6M
MTN (Ours)	MVDream	67.9	1.3 hours	8.3M

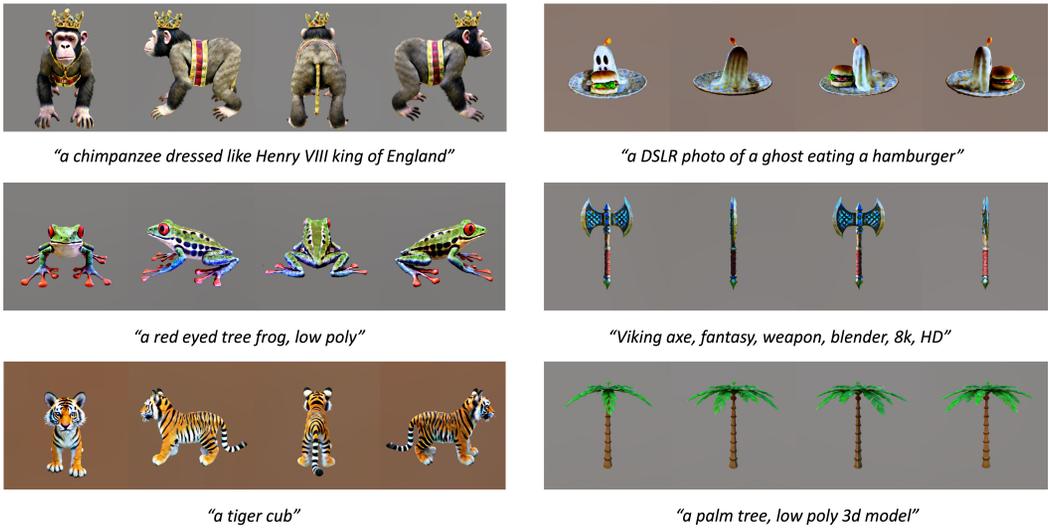


Fig. 10. Compatibility of the proposed method. Our method is highly compatible and can be easily scaled to other competitive multi-view diffusion models, such as MVDream [47], to further enhance the fidelity of 3D generation.

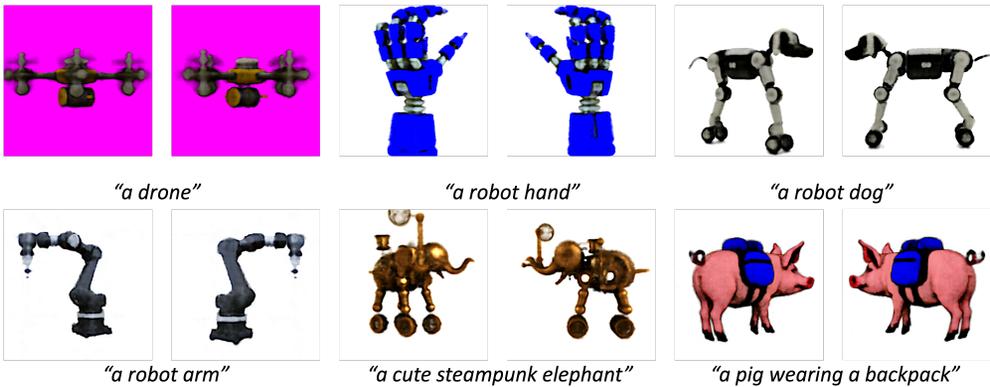


Fig. 11. Qualitative results on real-world and domain-specific prompts. Our model demonstrates strong generalization to challenging and abstract inputs (e.g., complex object compositions, artistic styles, and ambiguous language), producing geometrically consistent and visually realistic 3D structures.

Qualitative Evaluation on Real-world Prompts and Failure Cases. To further demonstrate adaptability and robustness in more realistic scenarios, we extend our evaluation to domain-specific prompts such as product design, architecture, and character modeling (see Figure 11). Our model maintains high geometric consistency and texture realism under these complex prompts, underscoring its generalization ability beyond simple synthetic cases. We also report several failure cases (Figure 12), where the generated results exhibit overly bright colors or excessive contrast. These issues likely arise from the diffusion supervision signal emphasizing high-frequency appearance cues during late-stage refinement. Mitigating this effect through tone-regularized guidance or exposure-balanced loss will be part of future work.

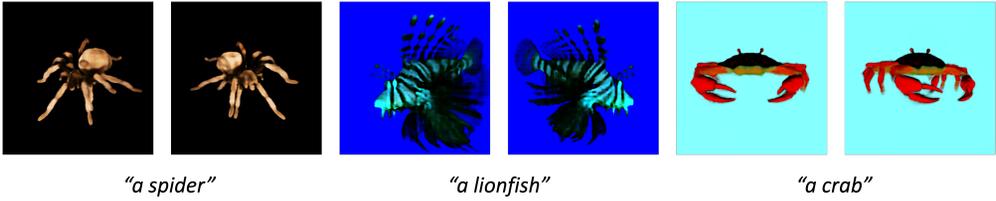


Fig. 12. Representative failure cases of our method. Some generated results exhibit excessive brightness, high contrast, or minor surface irregularities (e.g., color bleeding, topology distortions).

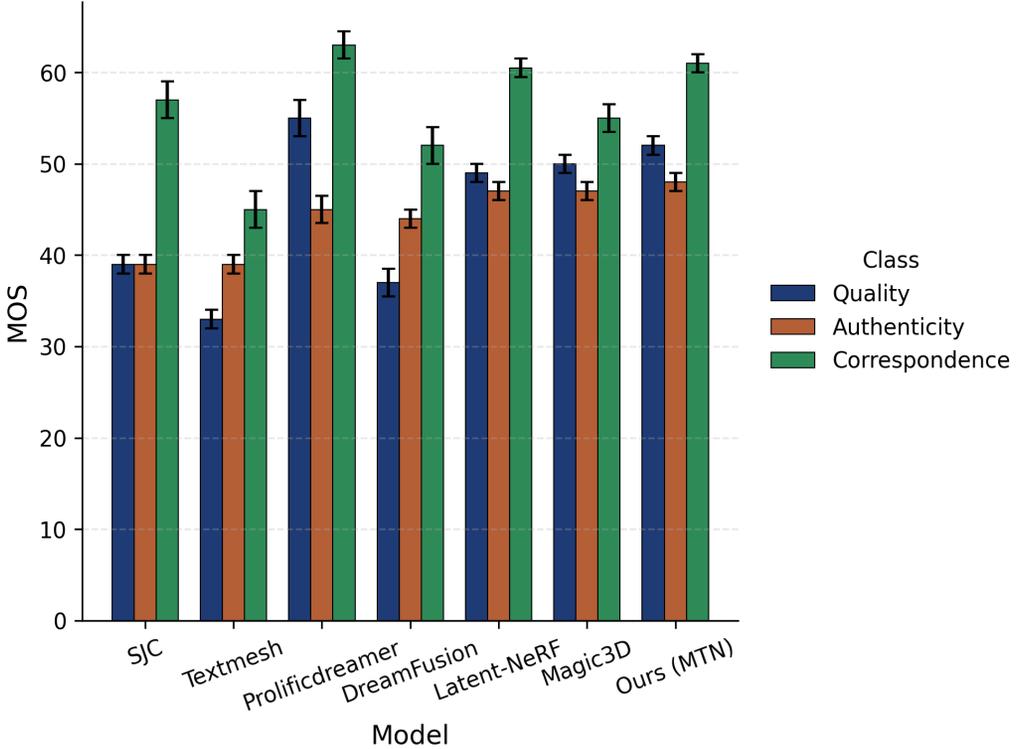


Fig. 13. Comparison of different text-to-3D generation methods evaluated on the T23DAQA benchmark across three dimensions: quality, authenticity, and correspondence.

Evaluation with 3D Quality Metrics. We further evaluate our generated 3D assets using the recently proposed T23DAQA [12], which benchmarks generation quality from three human-aligned dimensions: quality, authenticity, and text-asset correspondence. Following their evaluation protocol, we project each generated 3D asset into multi-view videos and compute model scores under all three dimensions. Under this metric, our method achieves the highest score on authenticity, indicating superior physical realism and plausible geometry (see Figure 13). Our quality and correspondence scores are also highly competitive, ranking only slightly below ProlificDreamer, while requiring significantly less training time. These results confirm that our model effectively balances efficiency and perceptual fidelity in text-to-3D generation.

Time cost and Model Size. Details are presented in Table 4. All experiments are performed on a V100 GPU. Among NeRF-based approaches, our model achieves the fastest convergence while

maintaining a significantly smaller parameter count. Although Gaussian Splatting-based methods train faster, they typically compromise generation quality and require additional post-processing to make the 3D outputs suitable for downstream applications such as 3D printing.

Table 4. Quantitative comparison with existing text-to-3D generation methods in terms of retrieval precision (R-Precision), parameter count, and training time. Our approach achieves the highest R-Precision among NeRF-based methods while using significantly fewer parameters and maintaining competitive efficiency.

Method	Type	R-Precision (%) \uparrow	#Params \downarrow	Training Time \downarrow
Latent-NeRF	NeRF-based	53.6	12.2 M	~ 1 hour
Stable-DreamFusion	NeRF-based	58.3	12.6 M	~ 1.5 hours
Magic3D	NeRF-based	62.0	13.2 M	2 hours
ProlificDreamer	NeRF-based	-*	15.1 M	> 20 hours
DreamGaussian	GS-based	61.6	-*	5 minutes
Ours	NeRF-based	63.3	8.3 M	~ 50 minutes

*: Due to the limitation of GPU resources, ProlificDreamer (153×20 hours) precision is unavailable, and GS-based models (e.g., DreamGaussian) are not directly comparable in parameter size due to architectural differences.

3D Printing. Our method provides a practical solution by directly converting the generated 3D output into a printable mesh format. The quality of these exported meshes is highlighted in Figure 1, which shows the uniformity of triangulation and the smoothness of surfaces. Such characteristics are crucial for the direct and efficient transmission of data to 3D printers. This is further evidenced by the six images on the right of Figure 14, showcasing the physical 3D products derived from these meshes. Our method significantly reduces the need for manual adjustments or additional post-processing steps, thereby streamlining the printing process. In addition, Figure 14 presents a detailed comparison of meshes rendered in Blender. Our method yields the most geometrically accurate and visually realistic models. In contrast, alternative methods exhibit noticeable distortions and structural inconsistencies, underscoring the robustness of our pipeline.

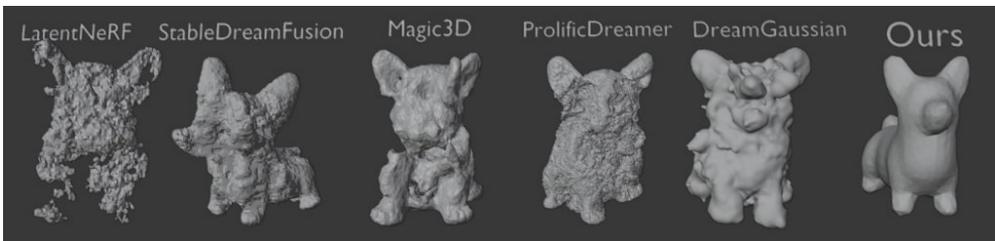


Fig. 14. Comparison of generated meshes (ready for print).

Potential downstream applications. Although our primary focus is on 3D printing, the exported textured meshes produced by our method are compatible with standard 3D editing and simulation pipelines, making them broadly applicable to other domains. For example, the reconstructed assets can be imported into Blender or Unity for AR/VR scene construction and game content design, where high-quality and editable 3D objects are essential. In robotics simulation, the physically consistent geometry enables integration into physics engines for task planning and interaction modeling. These extended applications further demonstrate the versatility and practical value of our framework beyond fabrication.

5 CONCLUSION

In this work, inspired by the bottom-up spirit, we introduce the Multi-scale Triplane Network (MTN) and a progressive learning strategy, both of which effectively ease the optimization difficulty during high-fidelity generation. The Multi-scale Triplane Network operates at the structure level to aggregate the multi-scale representation, while the progressive learning strategy functions at the recognition level to gradually refine high-frequency details. Extensive experiments verify the effectiveness of every component. We envision our approach offers a preliminary attempt for automatic 3D printing, bridging the gap between natural language descriptions and intricate 3D design. In the future, we will continue to explore the potential to complete occluded 3D objects [39] via language prior and discriminative language guidance [34].

REFERENCES

- [1] Mohammadreza Armandpour, Huangjie Zheng, Ali Sadeghian, Amir Sadeghian, and Mingyuan Zhou. 2023. Re-imagine the Negative Prompt Algorithm: Transform 2D Diffusion into 3D, alleviate Janus problem and Beyond. *arXiv preprint arXiv:2304.04968* (2023).
- [2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*. 41–48.
- [3] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. 2022. Efficient geometry-aware 3D generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16123–16133.
- [4] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. 2021. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5799–5809.
- [5] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. 2023. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 22246–22256.
- [6] Zilong Chen, Feng Wang, and Huaping Liu. 2024. Text-to-3d using gaussian splatting. In *CVPR*.
- [7] Yuhao Cheng, Yichao Yan, Wenhan Zhu, Ye Pan, Bowen Pan, and Xiaokang Yang. 2024. Head3d: Complete 3d head generation via tri-plane feature distillation. *ACM Transactions on Multimedia Computing, Communications and Applications* 20, 6 (2024), 1–20.
- [8] Zezhou Cheng, Menglei Chai, Jian Ren, Hsin-Ying Lee, Kyle Olszewski, Zeng Huang, Subhransu Maji, and Sergey Tulyakov. 2022. Cross-modal 3d shape generation and manipulation. In *European Conference on Computer Vision*. Springer, 303–321.
- [9] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems* 34 (2021), 8780–8794.
- [10] Lihe Ding, Shaocong Dong, Zhanpeng Huang, Zibin Wang, Yiyuan Zhang, Kaixiong Gong, Dan Xu, and Tianfan Xue. 2024. Text-to-3D Generation with Bidirectional Diffusion using both 2D and 3D priors. In *CVPR*.
- [11] Dylan Drotman, Saurabh Jadhav, Mahmood Karimi, Philip de Zonia, and Michael T Tolley. 2017. 3D printed soft actuators for a legged robot capable of navigating unstructured terrain. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 5532–5538.
- [12] Kang Fu, Huiyu Duan, Zicheng Zhang, Xiaohong Liu, Xionghuo Min, Jia Wang, and Guangtao Zhai. 2025. Multi-Dimensional Quality Assessment for Text-to-3D Assets: Dataset and Model. *arXiv preprint arXiv:2502.16915* (2025).
- [13] Tomoya Fujii, Jinqiang Dang, and Hiroto Tanaka. 2023. Hummingbird-bat hybrid wing by 3-D printing. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 3404–3410.
- [14] Lin Gao, Tong Wu, Yu-Jie Yuan, Ming-Xian Lin, Yu-Kun Lai, and Hao Zhang. 2021. Tm-net: Deep generative networks for textured meshes. *ACM Transactions on Graphics (TOG)* 40, 6 (2021), 1–15.
- [15] Lin Gao, Jie Yang, Tong Wu, Yu-Jie Yuan, Hongbo Fu, Yu-Kun Lai, and Hao Zhang. 2019. SDM-NET: Deep generative network for structured deformable mesh. *ACM Transactions on Graphics (TOG)* 38, 6 (2019), 1–15.
- [16] Kunal Gupta. 2020. *Neural mesh flow: 3d manifold mesh generation via diffeomorphic flows*. University of California, San Diego.
- [17] Paul Henderson, Vagia Tsiminaki, and Christoph H Lampert. 2020. Leveraging 2d data to learn textured 3d mesh generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 7498–7507.
- [18] Yukun Huang, Jianan Wang, Yukai Shi, Boshi Tang, Xianbiao Qi, and Lei Zhang. 2024. DreamTime: An Improved Optimization Strategy for Diffusion-Guided 3D Generation. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=1bAUywYJTU>

- [19] Moritz Ibing, Gregor Kobsik, and Leif Kobbelt. 2023. Octree transformer: Autoregressive 3d shape generation on hierarchically structured sequences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2697–2706.
- [20] Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. 2022. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 867–876.
- [21] Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4401–4410.
- [22] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics* 42, 4 (2023), 1–14.
- [23] Misha Konstantin. 2023. *DeepFloyd-IF*. <https://github.com/deep-floyd/IF>
- [24] Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. 2023. Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model. *arXiv preprint arXiv:2311.06214* (2023).
- [25] Jun Li, Kai Xu, Siddhartha Chaudhuri, Ersin Yumer, Hao Zhang, and Leonidas Guibas. 2017. Grass: Generative recursive autoencoders for shape structures. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 1–14.
- [26] Weiyu Li, Xuanyang Zhang, Zheng Sun, Di Qi, Hao Li, Wei Cheng, Weiwei Cai, Shihao Wu, Jiarui Liu, Zihao Wang, et al. 2025. Step1x-3d: Towards high-fidelity and controllable generation of textured 3d assets. *arXiv preprint arXiv:2505.07747* (2025).
- [27] Yixun Liang, Xin Yang, Jiantao Lin, Haodong Li, Xiaogang Xu, and Yingcong Chen. 2024. Luciddreamer: Towards high-fidelity text-to-3d generation via interval score matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6517–6526.
- [28] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. 2023. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 300–309.
- [29] Caixia Liu, Yali Chen, Minhong Zhu, Chenhui Hao, Haisheng Li, and Xiaochuan Wang. 2024. DEGAN: Detail-Enhanced Generative Adversarial Network for Monocular Depth-Based 3D Reconstruction. *ACM Transactions on Multimedia Computing, Communications and Applications* (2024).
- [30] Qihao Liu, Yi Zhang, Song Bai, Adam Kortylewski, and Alan Yuille. 2024. DIRECT-3D: Learning Direct Text-to-3D Generation on Massive Noisy 3D Data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6881–6891.
- [31] Ruoshi Liu, Rundui Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. 2023. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9298–9309.
- [32] Andrew Luo, Tianqin Li, Wen-Hao Zhang, and Tai Sing Lee. 2021. Surfgen: Adversarial 3d shape synthesis with explicit surface discriminators. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 16238–16248.
- [33] Yiwei Ma, Yijun Fan, Jiayi Ji, Haowei Wang, Haibing Yin, Xiaoshuai Sun, and Rongrong Ji. 2024. Creating High-quality 3D Content by Bridging the Gap Between Text-to-2D and Text-to-3D Generation. *ACM Transactions on Multimedia Computing, Communications and Applications* (2024).
- [34] Fumiya Matsuzawa, Yue Qiu, Kenji Iwata, Hirokatsu Kataoka, and Yutaka Satoh. 2023. Question Generation for Uncertainty Elimination in Referring Expressions in 3D Environments. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 6146–6152.
- [35] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. 2023. Latent-nerf for shape-guided generation of 3d shapes and textures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12663–12673.
- [36] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. 2022. Text2mesh: Text-driven neural stylization for meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13492–13502.
- [37] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.
- [38] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Tiberiu Popa. 2022. Clip-mesh: Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia 2022 conference papers*. 1–8.
- [39] Seyed S Mohammadi, Nuno F Duarte, Dimitrios Dimou, Yiming Wang, Matteo Taiana, Pietro Morerio, Atabak Dehban, Plinio Moreno, Alexandre Bernardino, Alessio Del Bue, et al. 2023. 3dsgrasp: 3d shape-completion for robotic grasp. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 3815–3822.
- [40] Charlie Nash, Yaroslav Ganin, SM Ali Eslami, and Peter Battaglia. 2020. Polygen: An autoregressive generative model of 3d meshes. In *International conference on machine learning*. PMLR, 7220–7229.

- [41] Alexander Quinn Nichol and Prafulla Dhariwal. 2021. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*. PMLR, 8162–8171.
- [42] Michael Niemeyer and Andreas Geiger. 2021. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11453–11464.
- [43] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. 2023. DreamFusion: Text-to-3D using 2D Diffusion. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=FjNys5c7VyY>
- [44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [45] Antoni Rosinol, Torsten Sattler, Marc Pollefeys, and Luca Carlone. 2019. Incremental visual-inertial 3d mesh generation with structural regularities. In *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 8220–8226.
- [46] Aditya Sanghi, Hang Chu, Joseph G Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshan. 2022. Clip-forge: Towards zero-shot text-to-shape generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18603–18613.
- [47] Yichun Shi, Peng Wang, Jianglong Ye, Long Mai, Kejie Li, and Xiao Yang. 2024. MVDream: Multi-view Diffusion for 3D Generation. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=FUgrjq2pbB>
- [48] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2021. Denoising Diffusion Implicit Models. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=St1giarCHLP>
- [49] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. 2020. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems* 33 (2020), 7537–7547.
- [50] Jiaxiang Tang. 2022. Stable-dreamfusion: Text-to-3D with Stable-diffusion. <https://github.com/ashawkey/stable-dreamfusion>.
- [51] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. 2024. DreamGaussian: Generative Gaussian Splatting for Efficient 3D Content Creation. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=UyNXMqnN3c>
- [52] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. 2017. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *Proceedings of the IEEE international conference on computer vision*. 2088–2096.
- [53] Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, Karsten Kreis, et al. 2022. LION: Latent Point Diffusion Models for 3D Shape Generation. *Advances in Neural Information Processing Systems* 35 (2022), 10021–10039.
- [54] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. 2022. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3835–3844.
- [55] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. 2023. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12619–12629.
- [56] Peng Wang and Yichun Shi. 2023. Imagedream: Image-prompt multi-view diffusion for 3d generation. *arXiv preprint arXiv:2312.02201* (2023).
- [57] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. 2023. ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [58] Bin-Shih Wu, Hong-En Chen, Sheng-Yu Huang, and Yu-Chiang Frank Wang. 2025. TPA3D: Triplane Attention for Fast Text-to-3D Generation. In *European Conference on Computer Vision*. Springer, 438–455.
- [59] Rundi Wu and Changxi Zheng. 2022. Learning to generate 3d shapes from a single example. *arXiv preprint arXiv:2208.02946* (2022).
- [60] Rundi Wu, Yixin Zhuang, Kai Xu, Hao Zhang, and Baoquan Chen. 2020. Pq-net: A generative part seq2seq network for 3d shapes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 829–838.
- [61] Xingyu Xie, Pan Zhou, Huan Li, Zhouchen Lin, and Shuicheng Yan. 2022. Adan: Adaptive nesterov momentum algorithm for faster optimizing deep models. *arXiv preprint arXiv:2208.06677* (2022).
- [62] Woo Yi Yang, Jiarui Wang, Sijing Wu, Huiyu Duan, Yuxin Zhu, Liu Yang, Kang Fu, Guangtao Zhai, and Xiongkuo Min. 2025. Lmme3dhf: Benchmarking and evaluating multimodal 3d human face generation with lmms. *arXiv preprint arXiv:2504.20466* (2025).
- [63] Taoran Yi, Jiemin Fang, Junjie Wang, Guanjun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and Xinggang Wang. 2024. GaussianDreamer: Fast Generation from Text to 3D Gaussians by Bridging 2D and 3D Diffusion Models. In *CVPR*.

- [64] Zibo Zhao, Zeqiang Lai, Qingxiang Lin, Yunfei Zhao, Haolin Liu, Shuhui Yang, Yifei Feng, Mingxin Yang, Sheng Zhang, Xianghui Yang, et al. 2025. Hunyuan3d 2.0: Scaling diffusion models for high resolution textured 3d assets generation. *arXiv preprint arXiv:2501.12202* (2025).
- [65] Xinyang Zheng, Yang Liu, Pengshuai Wang, and Xin Tong. 2022. SDF-StyleGAN: Implicit SDF-Based StyleGAN for 3D Shape Generation. In *Computer Graphics Forum*, Vol. 41. Wiley Online Library, 52–63.
- [66] Zhedong Zheng, Xiaohan Wang, Nenggan Zheng, and Yi Yang. 2022. Parameter-efficient person re-identification in the 3D space. *IEEE Transactions on Neural Networks and Learning Systems* 35, 6 (2022), 7534–7547. <https://doi.org/10.1109/TNNLS.2022.3214834> doi:10.1109/TNNLS.2022.3214834.
- [67] Zhedong Zheng, Jiayin Zhu, Wei Ji, Yi Yang, and Tat-Seng Chua. 2022. 3D Magic Mirror: Clothing Reconstruction from a Single Image via a Causal Perspective. *arXiv preprint arXiv:2204.13096* (2022).
- [68] Linqi Zhou, Yilun Du, and Jiajun Wu. 2021. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5826–5835.