# TextMesh: Generation of Realistic 3D Meshes From Text Prompts

Christina Tsalicoglou[1,2*]    Fabian Manhardt[2]    Alessio Tonioni[2]
Michael Niemeyer[2]    Federico Tombari[2,3]

[1]ETH Zurich    [2]Google    [3]Technical University of Munich

ctsalico@ethz.ch    {fabianmanhardt, alessiot, mniemeyer, tombari}@google.com

## Abstract

*The ability to generate highly realistic 2D images from mere text prompts has recently made huge progress in terms of speed and quality, thanks to the advent of image diffusion models. Naturally, the question arises if this can be also achieved in the generation of 3D content from such text prompts. To this end, a new line of methods recently emerged trying to harness diffusion models, trained on 2D images, for supervision of 3D model generation using view dependent prompts. While achieving impressive results, these methods, however, have two major drawbacks. First, rather than commonly used 3D meshes, they instead generate neural radiance fields (NeRFs), making them impractical for most real applications. Second, these approaches tend to produce over-saturated models, giving the output a cartoonish looking effect. Therefore, in this work we propose a novel method for generation of highly realistic-looking 3D meshes. To this end, we extend NeRF to employ an SDF backbone, leading to improved 3D mesh extraction. In addition, we propose a novel way to finetune the mesh texture, removing the effect of high saturation and improving the details of the output 3D mesh.*

## 1. Introduction

Generating photorealistic 2D images from simple text prompts is a rapidly growing field. Thanks to diffusion models and the availability of huge amount of training data with text-image pairs, current models can generate very high-quality images [25–27]. Naturally, the question arises if the same high quality generative capabilities can be achieved for 3D modeling. Unfortunately, this is a much more challenging field as the output space is significantly larger, 3D consistency is required, and there is a lack of large amount of training data pairs for text and 3D models.

Early methods mostly attempted at deforming template shapes, such as spheres, using a CLIP [24] objective. How-

Figure 1: **Exemplary results** of our TextMesh. Left: We compare our final mesh with the corresponding rendering from the public DreamFusion [23] gallery. While the results of DreamFusion are overly saturated, almost having a 'cartoonish' appearance, our mesh is more detailed and showcases a more realistic and natural appearance. Right: Since our method estimates a 3D mesh for the prompt instead of a NeRF-like representation, the obtained meshes can be directly plugged into standard computer graphics pipeline to *e.g.* enable AR/VR experiences.

ever, their emerging 3D shapes were still very unsatisfactory in geometry as well as appearance [10, 17, 19]. To overcome this limitation, DreamFusion [23] has recently proposed to harness the power of the aforementioned text-to-image diffusion models (*i.e.* Imagen [27]) to supervise 3D modelling from text prompts. To this end, they propose to

"a red-eyed tree frog"

"a plush triceratops toy, studio lighting, high resolution"

"a piglet sitting in a teacup"

"a pigeon reading a book"

"a swan and its cygnets swimming in a pond"

"a lemur taking notes in a journal"

"a pair of tan cowboy boots, studio lighting, product photography"

"a lion reading the newspaper"

"a tree stump with an axe buried in it"

"a teapot shaped like an elephant head where its snout acts as the spout"

Figure 2: **Qualitative Results.** Several qualitative 3D meshes generated from the given text prompts. The colors of the meshes are very natural, not showing any over-saturation effects.

train a Neural Radiance Field (NeRF) with a novel Score Distillation Sampling (SDS) gradient, together with view-dependant prompts. Despite their proposed method being capable of generating impressive results, it still has several downsides. First, the method has a tendency to produce objects with over-saturated colors due to the strong guidance required to make the model converge. Although prompt-engineering, *e.g.* prefixing "*A DSLR photo of [...]*" to the prompt, can mitigate this issue to some extent, the results are still not very satisfactory when it comes to actual realism. Second, [23] represents the 3D scene in the form of a NeRF, which renders the approach impractical to be used within standard computer graphics pipelines. Note that, while it is indeed possible to extract a mesh from NeRF, it is a non-trivial process given the density-based representation [31, 33].

In this work, we present TextMesh, a novel method for 3D shape generation from text prompts, targeted at tackling the aforementioned limitations, *i.e.* generating photorealistic 3D content in the form of standard 3D meshes. As demonstrated in Figure 1, our generated 3D meshes significantly improve upon [23] for realism and can be directly utilized within standard computer graphics pipelines and applications in AR or VR. To accomplish this, we modify DreamFusion to model radiance in the form of a signed distance function (SDF), allowing by design easy extraction of the surface as the 0-level set of the obtained volume. Furthermore, in an effort to enhance the mesh quality, we re-texture the output by leveraging another diffusion model,

conditioned on color and depth from the mesh. To this end, we render the object from multiple viewpoints and use diffusion to guide texture optimization to enhance realism and details. Nevertheless, when processing individual views independently, the refined texture exhibits severe inconsistencies. Therefore, we propose to run several views simultaneously through the diffusion model instead. To obtain the final texture, we then train on the produced output views together with Score Distillation Sampling to ensure smooth transitions. In Figure 2, we illustrate several meshes generated by our proposed method using different prompts.

To summarize, we propose the following contributions: **i)** We modify DreamFusion to model radiance in the form of SDF to tailor the model towards mesh extraction. **ii)** We propose a novel multi-view consistent and mesh conditioned re-texturing, enabling the generation of photorealistic 3D mesh models. **iii)** We experimentally show that our obtained meshes are geometrically of high quality and showcase more natural textures than the current state-of-the-art, whilst being ready to be deployed into pre-existing graphics pipelines.

## 2. Related work

**3D Reconstruction with Neural Fields.** Traditional 3D reconstruction methods [1, 4, 5, 11, 11, 30] usually rely on underlying depth [3, 28], or voxel [1, 4, 30]-based representations and perform some form of feature matching to fuse multi-view observations to a coherent 3D representation. While leading to satisfactory results in dense multi-

view stereo setups, these systems often fail in less constrained scenarios and cannot be integrated easily into other learning-based systems. In contrast, recent advances in neural fields have achieved impressive results on a variety of tasks. While seminal works focused on 3D reconstruction from 3D supervision [7, 16, 22], later works proposed surface rendering techniques [20, 34] that require 2D supervision in the form of image and mask data. The introduction of Neural Radiance Fields (NeRFs) [18] enabled impressive view synthesis from only image input via volume rendering. In many downstream applications, however, mesh-based representations are required, and directly extracting a mesh from a NeRF representation is non-trivial [31]. As a result, recent approaches [21, 32, 33] combine surface and volume rendering techniques to enable mesh extraction from image input. The goal of this work is to optimize a high-quality mesh and texture from text input. To this end, we adopt the VolSDF [33] representation due to its state-of-the-art performance and simple design.

**Photorealistic Image Generation From Text Prompts.**
Text-to-image models have recently achieved impressive high-fidelity and flexible image synthesis. The huge boost in quality has been made possible by the availability of extremely large datasets of image-text pairs [29] and scalable generator architectures based on diffusion models [2, 25–27] and transformers [6, 35]. One benefit of diffusion models over other classes of generative models is their flexibility with respect to the conditioning used in the image generation process since all of them support conditioning on text and a seed image. Recent works extend this support to text and a depth map [26] or text and a scene layout [8]. In this work, we utilize large-scale pretrained text-to-image models to enable text-to-3D mesh synthesis.

**3D Generation From Text Prompts.** There are a handful of works that attempt to generate 3D objects from text prompts. While most of them use a CLIP objective to supervise generation, a very new direction started to also incorporate large text-to-image diffusion models for training. As for the former, CLIPMesh [19] deforms a 3D sphere using a CLIP loss to obtain a 3D mesh that fits the input prompt, while Text2Mesh [17] similarly uses the CLIP loss to deform a given mesh and adjust its colors to better match the prompts. DreamFields [10] proposes to train a NeRF by rendering it from multiple viewpoints also using a CLIP objective. While these methods can indeed perform 3D object generation from text prompts, their results are very unsatisfactory when it comes to geometry and colors. Hence, more recently, DreamFusion started to investigate how to leverage pre-trained text-to-image diffusion models for 3D generation [23]. Similar to [10], DreamFusion trains a NeRF by rendering it from multiple viewpoints, however, supervis-

ing the model with their proposed Score Distillation Sampling (SDS) gradient based on *imagen*, a large text-to-image diffusion model [27]. While leading to impressive results, DreamFusion generates volumetric representations instead of meshes, making it impractical for many downstream applications such as graphics, where standard 3D representations such as meshes are required. Further, due to the high guidance weights required for optimization, their results tend to be oversaturated rather than photorealistic. To increase the texture resolution and improve the mesh extraction, the concurrent work Magic3D [14] proposes a two step approach. Firstly, they use a Dreamfusion like optimization with an SDF representation obtained from the density by subtracting a constant value. Secondly, they extract a mesh and employ differentiable rendering to further optimize it with a SDS objective.

In this work, we propose to overcome the aforementioned limitations by modifying the underlying neural field to represent an SDF instead of a radiance field to make the model better suited for mesh extraction and additionally propose a novel texture refinement to increase photorealism.

# 3. Method

Our goal is to develop a method that generates high-quality 3D mesh representations with photorealistic texture from text prompts. In the following, we discuss the main components of our method. We first discuss our initial neural field-based geometry and appearance representation (3.1) together with our first optimization stage where we train our model using a score-based distillation approach (3.2). Next, we describe how an initial mesh with texture can be extracted, and how we use this initial prediction to extract mesh-based RGB-D renderings. This information is then used in our second optimization stage, where we combine the output of an image and depth-conditioned diffusion model together with a score-based distillation approach to obtain our final mesh with photorealistic texture as output (3.3). Fig 3 presents an complete overview of our method.

## 3.1. Initial Scene Representation

**Neural Radiance Fields** A radiance field $f$ is a continuous mapping from a 3D location $\mathbf{x} \in \mathbb{R}^3$ and a ray viewing direction $\mathbf{d} \in \mathbb{S}^2$ to an RGB color $\mathbf{c} \in [0,1]^3$ and volume density $\sigma \in \mathbb{R}^+$. In Neural Radiance Fields (NeRF) [18], this field $f_\theta$ is parameterized as a neural network with parameters $\theta$. To render a pixel, a ray with direction $\mathbf{d} \in \mathbb{R}^2$ is cast from the camera center and $M$ equidistant points $\mathbf{p}_i$ are sampled along the ray. For a given camera pose $\xi \in \mathbb{R}^{3\times4}$, the operator $\pi$ maps a pixel $\mathbf{u} \in \mathbb{R}^2$ to its color $\hat{\mathbf{I}} \in [0,1]^3$
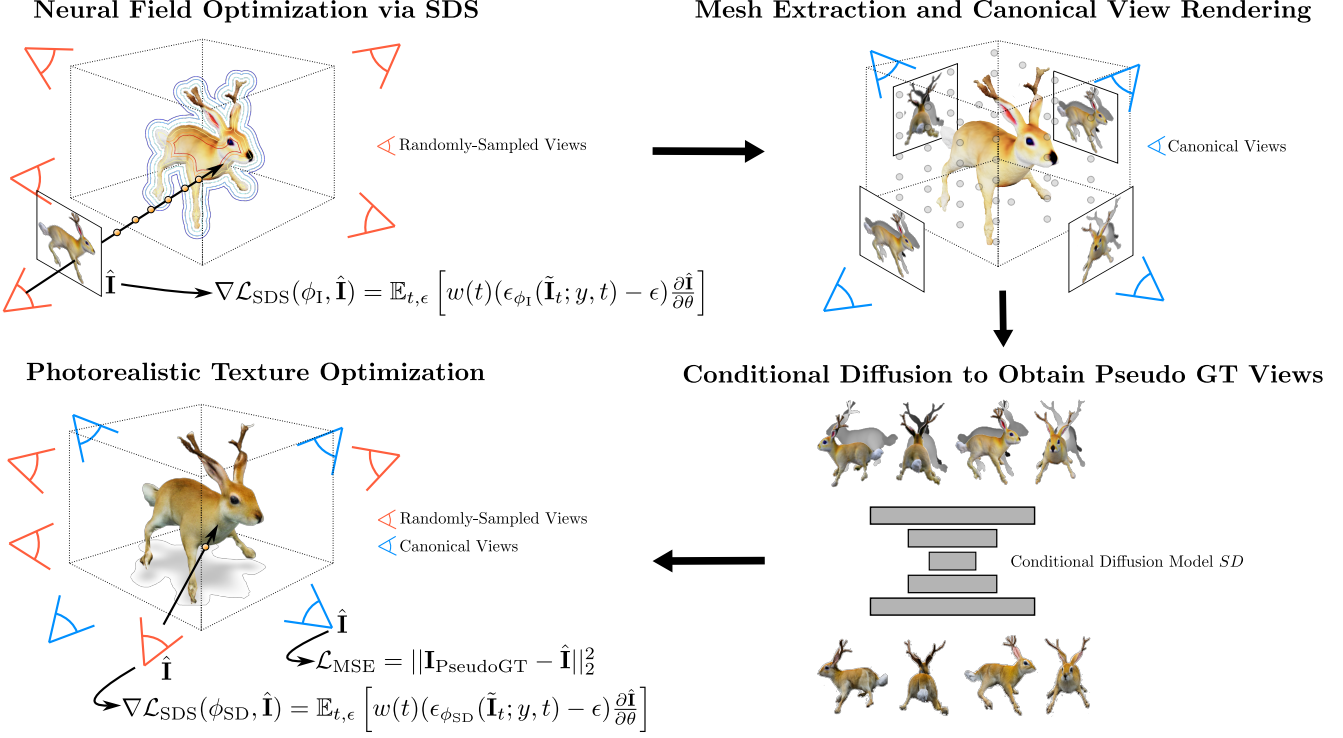
**Neural Field Optimization via SDS**

Randomly-Sampled Views

$$\nabla \mathcal{L}_{\text{SDS}}(\phi_{\text{I}}, \hat{\mathbf{I}}) = \mathbb{E}_{t,\epsilon}\left[w(t)(\epsilon_{\phi_{\text{I}}}(\tilde{\mathbf{I}}_t; y, t) - \epsilon)\frac{\partial \hat{\mathbf{I}}}{\partial \theta}\right]$$

**Mesh Extraction and Canonical View Rendering**

Canonical Views

**Photorealistic Texture Optimization**

Randomly-Sampled Views
Canonical Views

$$\mathcal{L}_{\text{MSE}} = ||\mathbf{I}_{\text{PseudoGT}} - \hat{\mathbf{I}}||_2^2$$

$$\nabla \mathcal{L}_{\text{SDS}}(\phi_{\text{SD}}, \hat{\mathbf{I}}) = \mathbb{E}_{t,\epsilon}\left[w(t)(\epsilon_{\phi_{\text{SD}}}(\tilde{\mathbf{I}}_t; y, t) - \epsilon)\frac{\partial \hat{\mathbf{I}}}{\partial \theta}\right]$$

**Conditional Diffusion to Obtain Pseudo GT Views**

Conditional Diffusion Model $SD$

Figure 3: **Schematic Overview.** Given the input text prompt *"an animal with the head of a rabbit, the body of a squirrel, the antlers of a deer, and legs of a pheasant"*, we train our initial distance field using Score Distillation Sampling (SDS) with view-dependent text prompting [23] and an Imagen prior (top left) and extract the mesh with marching cubes (top right). However, as the obtained appearance lacks details and the colors tend to be oversaturated, we render the color and depth from four orthogonal views of our mesh (top right) and run them jointly through StableDiffusion to generate photorealistic and 3D consistent views of our mesh (bottom right). Eventually, we finetune the mesh texture on the obtained views together with a small SDS gradient to account for minor misalignments (bottom left).

using classic volume rendering [18]:

$$\pi : (\xi, \mathbf{u}) \to \hat{\mathbf{I}}_u \;\; , \;\; \hat{\mathbf{I}}_u = \sum_{m=1}^{M} \alpha_m \mathbf{c}_m \tag{1}$$

where

$$\alpha_m = T_m \left(1 - \exp(-\sigma_m \delta_m)\right) \tag{2}$$

$$T_m = \exp\left(-\sum_{m'=1}^{m} \sigma_{m'} \delta_{m'}\right) \tag{3}$$

and $(\sigma_i, \mathbf{c}_i) = f_\theta(\mathbf{p}_i, \mathbf{d})$ are the evaluations along the ray and $\delta_i = ||\mathbf{p}_i - \mathbf{p}_j||_2$ are the Euclidean distances between sampled points.

**Signed Distance Fields** While NeRFs achieve impressive view synthesis results, the density-based representation is not well-suited for extracting a 3D geometry and obtaining a mesh [33]. To overcome this limitation, we instead adopt an SDF-based representation:

$$f_\theta(\mathbf{p}_i, \mathbf{d}) = (s_i, \mathbf{c}_i) \tag{4}$$

with $s_i \in \mathbb{R}$ being the signed distance from the surface at position $\mathbf{p}_i$. To enable training with volume rendering, we follow [33] and adopt the SDF to density transformation $t$:

$$t_\sigma(s) = \alpha \Psi_\beta(-s), \tag{5}$$

where

$$\Psi_\beta(s) = \begin{cases} \frac{1}{2}\exp\left(\frac{s}{\beta}\right) & \text{if } s \le 0 \\ 1 - \frac{1}{2}\exp\left(-\frac{s}{\beta}\right) & \text{if } s > 0 \end{cases} \tag{6}$$

with $\alpha, \beta \in \mathbb{R}$ being learnable parameters. Using this transformation, our SDF-based neural field representation can be rendered to the image plane using the same volume rendering technique from (1).

### 3.2. Text-to-3D via Score-based Distillation

We generate our initial 3D model via training a neural distance field using a score distillation sampling approach. To this end, given a randomly sampled camera pose $\xi$, we use our volume rendering operator $\pi$ from (1) on all pixels $\mathbf{u}_i$ on the image plane to obtain the respective rendered

image $\hat{\mathbf{I}}^2$. We then sample random normal noise and time step $t$ and add it to the rendered image using two weighting factors $\alpha_t$ and $\sigma_t$

$$\tilde{\mathbf{I}}_t = \alpha_t\hat{\mathbf{I}} + \sigma_t\epsilon \quad \text{where} \quad \epsilon \sim N(0, I) \tag{7}$$

Following [23], $\sigma_t$ is chosen such that $\tilde{\mathbf{I}}_t$ is close to the data density at the start of the diffusion process, i.e., $\sigma_0 \approx 0$ and converging to 1 for maximum diffusion steps, while $\alpha_t^2 = 1 - \sigma_t^2$. We feed $\tilde{\mathbf{I}}$ to a diffusion model $\phi_{\mathrm{I}}$ (Imagen [27] in our experiments), which attempts to predict the noise $\epsilon$ with $\epsilon_{\phi_{\mathrm{I}}}(\tilde{\mathbf{I}}, y, t)$, given the noisy image $\tilde{\mathbf{I}}$, diffusion step $t$, and text embedding $y$. From this prediction we can derive the gradient direction pushing the rendered images to a high probability density region for the provided text prompt with

$$\nabla\mathcal{L}_{\mathrm{SDS}}(\phi_{\mathrm{I}}, \hat{\mathbf{I}}) =$$
$$\mathbb{E}_{t,\epsilon}\left[w(t)(\epsilon_{\phi_{\mathrm{I}}}(\tilde{\mathbf{I}}_t; y, t) - \epsilon)\frac{\partial\hat{\mathbf{I}}}{\partial\theta}\right], \tag{8}$$

where $w(t)$ denotes a weighting function and $y$ is the conditioning text embedding. This gradient is then utilized to optimize our signed distance field till convergence. Similar to [14] and [23], we also employ classifier-free guidance [9] to control the strength of the text conditioning. Since this process involves rendering a MLP-based NeRF volume and running a pixel level diffusion model it can be carried out only at low resolution due to memory constraint. For this reason we compute $\mathcal{L}_{\mathrm{SDS}}$ on images rendered at $64\times64$ (i.e., we use only the low resolution branch of Imagen). Note that the exact details for rendering, including shading and background modeling, match those from Dream-Fusion [23], which we omitted for the sake of clarity. We kindly refer to their paper for more information. However, unlike [23], we sample the whole elevation range for the camera, to avoid bleeding artifacts at the model bottom.

Eventually, a mesh is extracted from the signed distance field as the surface at the zero-level set using Marching Cubes (MC) [15]. Since floaters (i.e. areas of near 0 signed distance value away from the expected object surface) can occasionally remain within the volume, we additionally always select the largest mesh component closer to the center of the volume to create a mesh and to be used for the following steps.

### 3.3. Photorealistic Texturing Using Multi-View Consistent Diffusion

Upon having extracted the 3D mesh $\mathcal{M}$ from our trained distance field we already have a good geometry for the model. On the other hand, the texture still includes two main drawbacks: it misses high frequency details since the optimization in 3.2 is performed at low resolution only,

(a) Input: Mesh rendering.



(b) Output: Processed Independently.
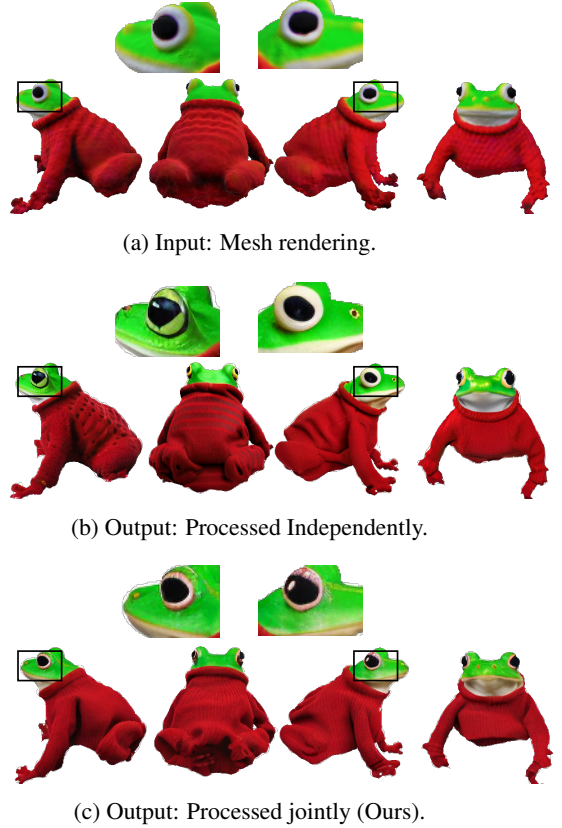


(c) Output: Processed jointly (Ours).

Figure 4: **Diffusion model conditioning strategies**. Input (a) and outputs of the depth conditioned diffusion model processing each input image independently (b) and conditioning on all four input images jointly (c). We observe that the latter leads to multi-view consistent results, while processing each image independently can introduce inconsistencies (see zoom boxes).

and it shows over-saturated ('cartoonish'), colors as the result of using a large guidance weight [9]. To solve these two limitations we refine the initial texture using the standard pipeline of a Stable Diffusion model $SD$ [26] conditioned on color and depth. To this end, we take our obtained mesh, freeze its geometry, and use a differentiable-render $\mathcal{R}$ (NVdiffrast [13]) to render color and depth from four canonical viewpoints $\mathcal{P}$ (i.e., front, back, and both sides). Feeding the four views independently to a depth-conditioned diffusion model would be a straightforward way to obtain highly realistic images of the object, which could serve to guide the re-texturing. However, when processed independently, the resulting images exhibit several 3D inconsistencies, which would give the object a different identity depending on the viewpoint.

To overcome this limitation, we propose to tile the four canonical RGB and depth predictions on a $2 \times 2$ grid to a

single RGB image $\hat{\mathbf{I}}_{\text{tiled}}$ and depth map $\mathbf{D}_{\text{tiled}}$ and process them jointly in a single diffusion operation

$$\mathbf{I}_{\text{tiled}} = SD(\hat{\mathbf{I}}_{\text{tiled}}, \mathbf{D}_{\text{tiled}}). \quad (9)$$

This enforces the diffusion model to generate consistent views during the diffusion of the tiled image (See Fig. 4). The individual pseudo GT views $\{\mathbf{I}_{\text{PseudoGT}, i}\}_{i=1}^{4}$ are extracted from image $\mathbf{I}_{\text{tiled}}$ and then serve as a pseudo ground truth that allows us to apply the new texture to the mesh geometry. The loss we optimize is

$$\mathcal{L}_{\text{texture}}(\mathcal{R}, \mathcal{M}, P, i) = ||\mathbf{I}_{\text{PseudoGT},i} - \hat{\mathbf{I}}||_2^2 \quad (10)$$

$$\text{with} \quad \hat{\mathbf{I}} = \mathcal{R}(\mathcal{M}, P) \quad \text{for } P \in \mathcal{P} \quad (11)$$

While tiling the views significantly improves 3D object consistency, the views can still exhibit minor misalignment at their intersection as well as on unobserved object parts. To ensure smooth transitions and a complete 3D mesh, we perform a second optimization stage, where we combine a photometric loss with a small SDS component using an image-to-image Stable Diffusion model $\phi_{\text{SD}}$. The new pseudo ground truth for this stage, $\{\mathbf{I}'_{\text{PseudoGT},i}\}_i$, is obtained by rendering the converged texture from poses $\mathcal{P}'$. In this stage, we then optimize

$$\nabla\mathcal{L}_{\text{texture}}(\mathcal{R}, \mathcal{M}, P, i) = \nabla\mathcal{L}_{\text{MSE}} + \lambda_{\text{SDS}}\nabla\mathcal{L}_{\text{SDS}}, \quad (12)$$

$$\text{with} \quad \mathcal{L}_{\text{MSE}} = ||\mathbf{I}'_{\text{PseudoGT},i} - \hat{\mathbf{I}}||_2^2 \quad (13)$$

$$\text{and} \quad \hat{\mathbf{I}} = \mathcal{R}(\mathcal{M}, P) \quad \text{for } P \in \mathcal{P}' \quad (14)$$

where $i$ describes the viewpoint for the camera poses $P \in P'$ and and $\lambda_{\text{SDS}}$ is the SDS weighting that controls its contribution to the texture optimization. For this step, we use a very small guidance weight of 7.5 compared to previous only SDS-supervised works, as it has been reported that increasing the guidance weight often results in saturated colors [9] and we only want to make small changes to the texture. Further, anchoring the optimization on $\mathbf{I}'_{\text{PseudoGT},i}$ enforces the resulting texture to not deviate too much from the original, encouraging only regions with high SDS gradients to change.

## 4. Evaluation

### 4.1. Experimental Setup

**Metrics** We follow previous work [10, 19, 23] and evaluate our method using the CLIP [24] R-Precision metric. This metric measures how well images rendered from the generated geometry correlate with the provided input text prompt; however, it fails to capture any aspect related to the 3D consistency and photorealistic appearance of the generated shapes. Therefore, we additionally report the $\text{FID}_{\text{CLIP}}$ score [12], which evaluates the FID in the feature space of

| Method | CLIP R-Precision ↑ | | | $\text{FID}_{\text{CLIP}}$ ↓ |
| --- | --- | --- | --- | --- |
| | B/32 | B/16 | L/14 | |
| CLIP-Mesh | 100 | 100 | 99.0 | 57.5 |
| DreamFusion | 94.3 | 97.1 | 97.1 | 59.3 |
| Ours | 91.4 | 91.4 | 94.3 | 57.4 |

Table 1: **Comparison with state of the art.** Comparing our method against the state of the art using R-Precision for different CLIP models and $\text{FID}_{\text{CLIP}}$. The metrics are computed on 35 prompts from the public DreamFusion gallery.

the CLIP ViT-B-32 image encoder. As reference images, we use the ImageNet 2012 validation set. Similar to DreamFusion, we render 60 azimuthal angles at an elevation of 30 degrees to generate images which are used to evaluate both the R-Precision and $\text{FID}_{\text{CLIP}}$.

**Text prompts** We use 35 text prompts available on the public DreamFusion gallery[3], which mainly contains prompts describing individual objects making them suitable to be turned into meshes.

### 4.2. Comparison with state-of-the-art

We first compare our approach quantitatively against state-of-the-art methods in Table 1. We compare to CLIP-Mesh as an alternative text-to-mesh generation model and to DreamFusion as a state-of-the-art text-to-NeRF method. For this comparison we re-run all the competitors on the same 35 prompts using the original code kindly provided by the authors and the default settings. Interestingly, although producing the qualitative worst results, CLIP-Mesh is able to report the best quantitative numbers for R-Precision. This is explained by the fact that CLIP-Mesh directly optimizes for the CLIP metric during its training. Our method performs on par with Dreamfusion, obtaining somewhat worse R-precision and better $\text{FID}_{\text{CLIP}}$ results.

We also provide a qualitative comparison to the competitors in Figure 5. For this comparison we also include Magic3D [14], taking their results from their paper. The qualitative comparison indicates that, while all methods allow for some degree of text-to-3D generation, our method achieves more photorealistic results. While all baselines tend to have a cartoonish and oversaturated appearance, we achieve more natural texture predictions thanks to our proposed texturing stage. Note that a photorealistic appearance is crucial for many AR/VR applications where 3D content should fit in smoothly with the environment.
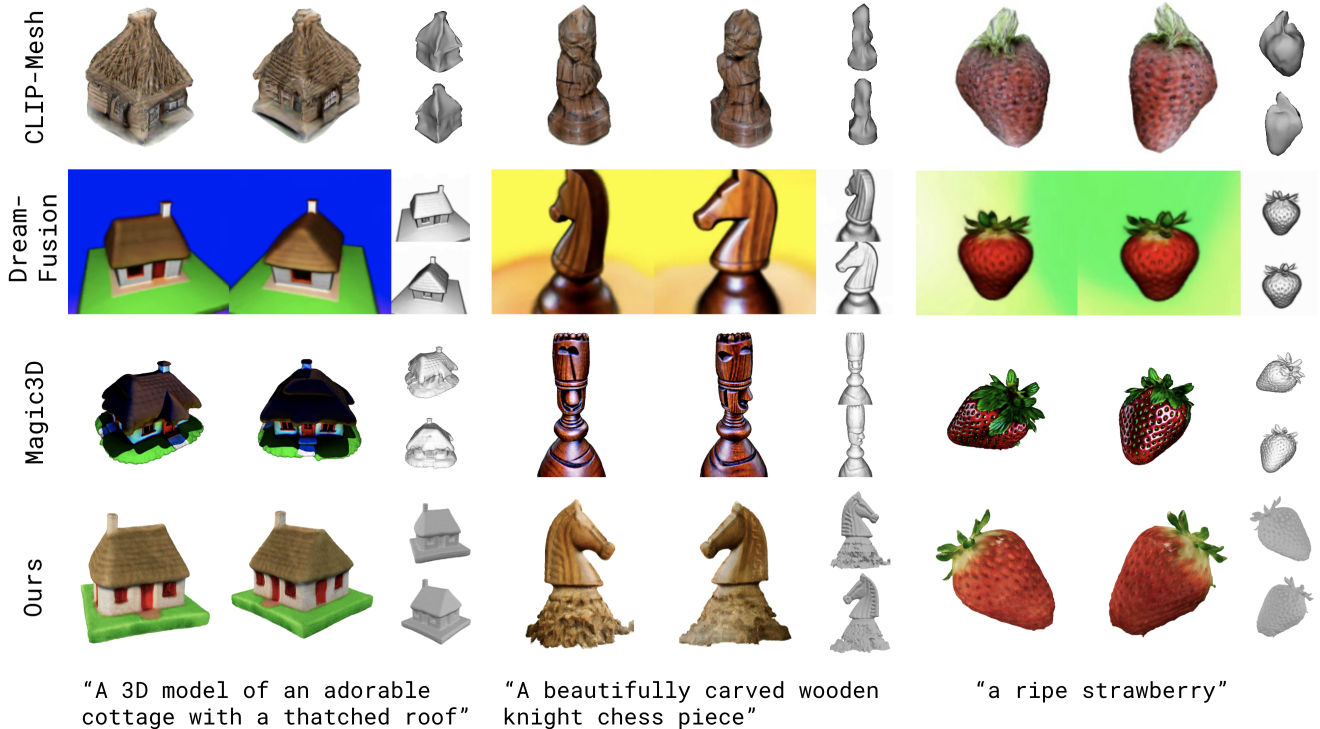
Figure 5: **Qualitative comparison** between our meshes, the meshes by CLIP-Mesh [19] and Magic3d [14], and the NeRF volumes by DreamFusion [23]. For each model we show both RGB renderings and 3D shape. Results for DreamFusion and Magic3D from [14].

| Method | R-Precision ↑ CLIP B/32 | FID$_{\text{CLIP}}$ ↓ . |
|---|---|---|
| w/o texture finetuning | 85.7 | 61.1 |
| w/o depth conditioning | 91.4 | 57.7 |
| w/o joint diffusion | 91.4 | 55.9 |
| w/o multi-view loss | 80.0 | 61.1 |
| Ours | 91.4 | 57.4 |

Table 2: **Ablation** of various components of our method using R-Precision and FID$_{\text{CLIP}}$ on 35 meshes.

| Criteria | Preference (%) ↑ |
|---|---|
| More Natural Colors | 61.2 |
| More Detailed Texture | 63.3 |
| Overall Visually Preferred | 57.9 |

Table 3: **User Study.** Results of our user study, conducted with 30 participants. Each participant was shown two meshes, before and after our re-texturing using depth-condition diffusion, for a total of 15 prompts from the DreamFusion [23] gallery and had to choose which mesh they preferred based on different criteria.

## 4.3. Ablation Study

In Table 2 we ablate various components of our method. Each ablation is evaluated on 35 prompts and, as before, we report the CLIP R-Precision and FID$_{\text{CLIP}}$ score. In our default setting, we use a depth-conditioned Stable Diffusion model and four views of the mesh seen from the front, back and sides (see Figure 3). The views are then tiled in a grid and processed as one image by Stable Diffusion.

**Quantitative Results** First, the results in Table 2 indicate that the texture finetuning stage is crucial to obtain realistic-looking meshes, as all refined options obtain better R-Precision and FID$_{\text{CLIP}}$ scores. Further, when removing depth conditioning or joint diffusion, we obtain overall similar results, with separate diffusion obtaining slightly better results for FID$_{\text{CLIP}}$. We attribute this to the fact that the employed metrics, including FID$_{\text{CLIP}}$, are not very suitable at evaluating the 3D consistency of the generated texture. We kindly refer to the supplement, where we provide several examples demonstrating that removing joint diffusion

---

[3]https://dreamfusion3d.github.io/gallery.html

Figure 6: **3D Consistency.** Our depth-conditioned joint diffusion process ensures that obtained meshes are geometrically accurate and consistent in 3D space.



Figure 7: **Comparing the 3D mesh geometry** from the radiance field of DreamFusion and our SDF-based approach.

leads to inconsistencies in 3D space. Finally, the multi-view component is essential for obtaining realistic results, as the SDS-only driven optimization performs worst overall.

**User Study**   While the reported metrics can provide an indication of the quality of the results, it is important to note that they do not directly measure the perceived quality of the generated models. To further quantify the importance of our texture finetuning stage, we perform a user study comparing the results before and after the photorealistic texturing stage in Table 3. We observe that humans prefer the results after the texturing stage, in particular with respect to texture details and color.

## 4.4. Mesh Quality

**Extracting Meshes**   In Figure 7, we compare extracted meshes from Dreamfusion [23] and our method. We observe that our SDF-based approach leads to smoother mesh predictions. Obtaining a high-quality mesh as output representation is crucial for many applications, and we provide an example for an AR application in Figure 1. To generate the meshes for DreamFusion we use marching cubes over the volumes obtained for the evaluations in Table 1, i.e., applying the default settings provided by the authors. Besides the use of a SDF volume, a key difference between our method and Dreamfusion is that we sample the full elevation range of camera poses while DreamFusion uses a limited one. This allows us to obtain nice complete meshes without spurious surfaces (e.g., the artefacts on the bottom of the DreamFusion lion).

**3D Consistency**   In Figure 6, we show multiple views for the same object to provide a qualitative evaluation of the 3D consistency of the optimized meshes with photorealistic texture. We find that the final outputs of our method are indeed 3D consistent and that the texture appears realistic from arbitrary viewpoints.

## 5. Conclusion

We present TextMesh, a novel approach for 3D mesh generation from text prompts. In the core, we propose to represent the geometry as a distance field which is optimized using Score Distillation Sampling (SDS). After optimization of the distance field, we then extract the mesh and refine its original texture to achieve a more detailed and natural appearance. In contrast to similar methods, we supervise the texture refinement primarily with a photometric

loss on enhanced 2D mesh renderings generated by a depth condition image-to-image diffusion model, and only rely on SDS to smooth out transitions within our multi-view supervision. This leads to more photorealistic textures, preferred by a larger portion of our survey participants compared to the initial, unrefined texture.

# References

[1] Motilal Agrawal and Larry S Davis. A probabilistic framework for surface reconstruction from multiple images. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2001. 2

[2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 3

[3] Michael Bleyer, Christoph Rhemann, and Carsten Rother. Patchmatch stereo - stereo matching with slanted support windows. In *Proc. of the British Machine Vision Conf. (BMVC)*, 2011. 2

[4] JS De Bonet and Paul Viola. Poxels: Probabilistic voxelized volume reconstruction. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 1999. 2

[5] A Broadhurst, T W Drummond, and R Cipolla. A probabilistic framework for space carving. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2001. 2

[6] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. 3

[7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3

[8] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pages 89–106. Springer, 2022. 3

[9] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 5, 6

[10] 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. 1, 3, 6

[11] Kiriakos N. Kutulakos and Steven M. Seitz. A theory of shape by space carving. *International Journal of Computer Vision (IJCV)*, 38(3):199–218, 2000. 2

[12] Tuomas Kynkäänniemi, Tero Karras, Miika Aittala, Timo Aila, and Jaakko Lehtinen. The role of imagenet classes in fr\'echet inception distance. *arXiv preprint arXiv:2203.06026*, 2022. 6

[13] Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. Modular primitives for high-performance differentiable rendering. *ACM Transactions on Graphics (TOG)*, 39(6):1–14, 2020. 5

[14] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. *arXiv preprint arXiv:2211.10440*, 2022. 3, 5, 6, 7

[15] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987. 5

[16] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3

[17] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13492–13502, 2022. 1, 3

[18] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2020. 3, 4

[19] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Tiberiu Popa. Clip-mesh: Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–8, 2022. 1, 3, 6, 7

[20] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3

[21] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 3

[22] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3

[23] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 1, 2, 3, 4, 5, 6, 7, 8

[24] 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. 1, 6

[25] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1, 3

[26] 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. 1, 3, 5

[27] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 1, 3, 5

[28] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2016. 2

[29] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022. 3

[30] S.M. Seitz and C.R. Dyer. Photorealistic scene reconstruction by voxel coloring. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 1997. 2

[31] Jiaxiang Tang, Hang Zhou, Xiaokang Chen, Tianshu Hu, Errui Ding, Jingdong Wang, and Gang Zeng. Delicate textured mesh recovery from nerf via adaptive surface refinement. *arXiv preprint arXiv:2303.02091*, 2022. 2, 3

[32] Shaofei Wang, Marko Mihajlovic, Qianli Ma, Andreas Geiger, and Siyu Tang. Metaavatar: Learning animatable clothed human models from few depth images. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 3

[33] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. 2, 3, 4

[34] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. In *Advances in Neural Information Processing Systems (NIPS)*, 2020. 3

[35] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 3