# Dream3D: Zero-Shot Text-to-3D Synthesis Using 3D Shape Prior and Text-to-Image Diffusion Models

Jiale Xu[1,3*]    Xintao Wang[1†]    Weihao Cheng[1]    Yan-Pei Cao[1]

Ying Shan[1]    Xiaohu Qie[2]    Shenghua Gao[3,4,5†]

[1]ARC Lab, [2]Tencent PCG    [3]ShanghaiTech University

[4]Shanghai Engineering Research Center of Intelligent Vision and Imaging

[5]Shanghai Engineering Research Center of Energy Efficient and Custom AI IC

Figure 1. Thanks to the incorporation of 3D shape prior and powerful text-to-image diffusion models, our method Dream3D is capable of generating 3D content with better visual quality and shape accuracy conforming to text prompt than PureCLIPNeRF [31].

## Abstract

*Recent CLIP-guided 3D optimization methods e.g., DreamFields [25] and PureCLIPNeRF [31] achieve great success in zero-shot text-guided 3D synthesis. However, due to the scratch training and random initialization without any prior knowledge, these methods usually fail to generate accurate and faithful 3D structures that conform to the corresponding text. In this paper, we make the first attempt to introduce the explicit 3D shape prior to CLIP-guided 3D optimization methods. Specifically, we first generate a high-quality 3D shape from input texts in the text-to-shape stage as the 3D shape prior. We then utilize it as the initialization of a neural radiance field and then optimize it with the full prompt. For the text-to-shape generation, we present a simple yet effective approach that directly bridges the text and image modalities with a powerful text-to-image diffusion model. To narrow the style domain gap between images synthesized by the text-to-image model and shape renderings used to train the image-to-shape generator, we further propose to jointly optimize a learnable text prompt and fine-tune the text-to-image diffusion model for rendering-style image generation. Our method, namely, Dream3D, is capable of generating imaginative 3D content with better visual quality and shape accuracy than state-of-the-art methods. Our Project page is at https://bluestyle97.github.io/dream3d/.*

## 1. Introduction

Text-guided 3D synthesis aims to generate 3D content consistent with an input text, which has the potential to benefit a broad range of applications, *e.g.*, animations, games, and virtual reality. Modern zero-shot text-to-image models [11, 40, 51, 52, 55, 58] have achieved unprecedented success and can synthesize diverse, high-fidelity, and imaginative images from various text prompts. However, it is challenging to replicate the growth in text-to-image synthesis to the text-to-3D synthesis task since collecting a large-scale paired text-3D dataset is unfeasible.

Zero-shot text-guided 3D synthesis [25, 28, 31, 33, 47, 59], on the contrary, usually leverages powerful vision-language models, *e.g.*, CLIP [69], to eliminate the demand for paired data, and thus is more attractive. Typically, there are

---

*Work done during an internship at ARC Lab, Tencent PCG.

†Corresponding Author.

1

two categories. **1)** CLIP-based generative models [33, 59]. CLIP-Forge [59] is a representative method, which uses images as an intermediate bridge. It trains a mapper from the CLIP image embeddings of ShapeNet renderings to the shape embeddings of a 3D shape generator, then switches to CLIP text embedding as the mapper's input at test time. **2)** CLIP-guided 3D optimization methods, *e.g.*, Dream-Fields [25] and PureCLIPNeRF [31]. They continuously optimize the CLIP similarity loss between the text prompt and rendered images of a 3D scene representation, *e.g.*, neural radiance fields [38]. The first category heavily relies on 3D shape generators trained on limited 3D shapes and seldom has the capacity to adjust its shape structures, while the second category has more creative freedom with the "dreaming ability" to generate diverse shape structures and textures.

We start to develop our method from CLIP-guided 3D optimization methods. Though these methods can produce impressive results, we observe that they usually fail to generate accurate and faithful 3D structures that conform to the corresponding texts (Fig. 1, $2^{nd}$ row). Due to the scratch training and random initialization without any prior knowledge, these methods tend to generate highly-unconstrained "adversarial content" that have high CLIP scores but low visual quality. To overcome the aforementioned issue and synthesize more faithful 3D content, we propose to first generate a high-quality 3D shape from input texts, and then leverage it as an explicit **"3D shape prior"** to the CLIP-guided 3D optimization process. Specifically, in the *text-to-shape** stage, we first synthesize a 3D shape without textures of the main common object in the text prompt. We then utilize it as the *initialization* of a voxel-based neural radiance field and then optimize with the full prompt.

The *text-to-shape* generation itself is a challenging task. Previous methods [25, 59] are often trained on images and tested with texts, and adopt CLIP to bridge the two modalities. Such a design leads to a mismatching problem due to the gap between the CLIP text and image embedding spaces. Moreover, existing methods are not capable of generating high-quality 3D shapes. In this work, we propose to directly bridge the text and image modalities with a powerful text-to-image diffusion model, *e.g.*, Stable Diffusion [55]. Specifically, we adopt the *text-to-image* diffusion models to synthesize an image from the input text, and then feed the image into an *image-to-shape* generator [79] to produce high-quality 3D shapes. Due to the same procedure both in training and testing, the mismatching problem can be largely reduced. However, it still suffers from the style domain gap between the images synthesized by the text-to-image model and the shape renderings used to train the

image-to-shape generator. Inspired by recent work on controllable text-to-image synthesis [14, 57], we propose to address this domain gap by jointly optimizing a learnable text prompt and fine-tuning Stable Diffusion to obtain a stylized prompt and a stylized generator. Thus, we can synthesize images that meet the style of shape renderings we used to train the image-to-shape module, without suffering from the domain gap.

To summarize, **1)** We make the first attempt to introduce the explicit **3D shape prior** into CLIP-guided 3D optimization methods. The proposed method can generate more accurate and high-quality **3D shapes** conforming to the corresponding text, while still enjoying the **"dreaming" ability** of generating diverse shape structures and textures (Fig. 1, $1^{st}$ row). Therefore, we name our method "Dream3D" as it has both two strengths. **2)** For text-to-shape generation, we present a simple yet effective approach that directly bridges the text and image modalities with a powerful text-to-image diffusion model. In order to narrow the style domain gap between generated images and shape renderings, we further jointly optimize a learnable text prompt and fine-tune a text-to-image diffusion model for rendering-style image generation. **3)** Our Dream3D is capable of generating imaginative 3D content with better visual quality and shape accuracy than state-of-the-art methods. Moreover, our text-to-shape pipeline can also generate 3D shapes with higher quality than previous work.

## 2. Related Work

**3D Shape Generation.** Generative models on 3D shapes have been extensively studied in recent years. It is more challenging than 2D image generation due to the expensive 3D data collection and the complexity of 3D shapes. Different 3D generative models adopt varied shape representations, *e.g.*, voxel grids [32, 67], point clouds [1, 36, 48, 75, 77, 80], meshes [16, 17, 20, 21, 39], and implicit fields [7, 8, 30, 35, 70, 71, 79]. They are trained to model the distribution of shape geometry (and optionally, texture) from a collection of 3D shapes. In comparison, some methods [2, 3, 10, 15, 19, 41, 43, 45, 46, 62, 73, 74] seek to relax the demand for direct 3D supervision and learn a 3D generator from 2D image collections. These methods incorporate explicit 3D representations, *e.g.*, meshes [15, 45, 46] and neural radiance fields [2, 3, 10, 19, 38, 41, 43, 62, 73, 74], with surface or volume based differentiable rendering techniques [6, 26, 27, 38, 42] to enabling the learning of 3D awareness from images.

**Text-to-Image.** Previous attempts on text-to-image synthesis [49, 53, 56, 66, 72, 78] mainly focus on some domain-specific datasets, and usually adopt the Generative Neural Networks (GANs) [18] architecture. Benefiting from large-scale paired image-text datasets [60, 61] and scalable generative architectures, zero-shot text-to-image synthesis has

---

*We use the term "shape" to denote the 3D geometric models *without textures* throughout the paper, while some works [5, 34] also use this term for textured 3D models.
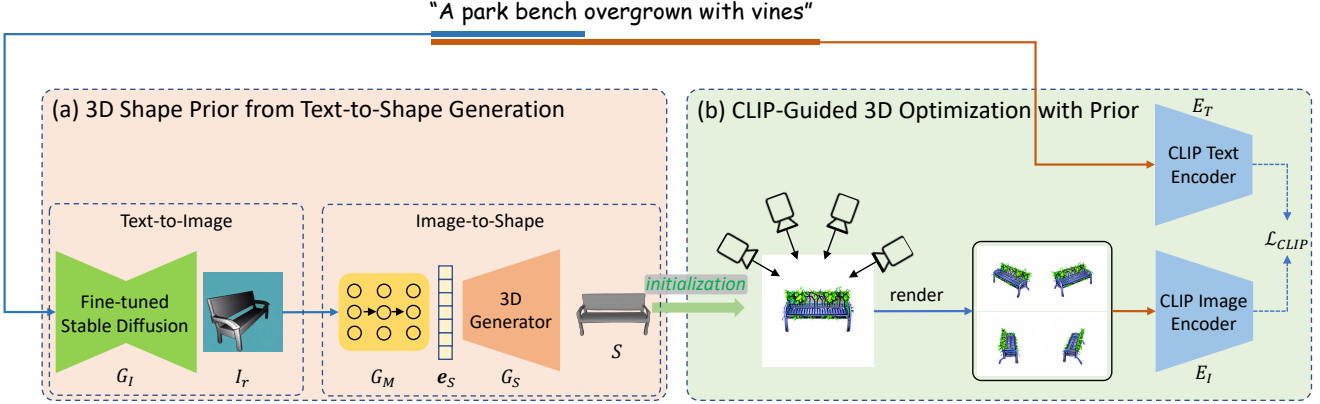
"A park bench overgrown with vines"

Figure 2. Overview of our text-guided 3D synthesis framework. (a) In the first text-to-shape stage (Sec. E), we leverage a fine-tuned off-the-shelf text-to-image diffusion model $G_I$ to synthesize a rendering-style image $I_r$ from the input text prompt $y$. Then we generate a latent shape embedding $e_S$ from $I_r$ using a diffusion-model-based shape embedding generation network $G_M$. Finally, we feed $e_S$ into a high-quality 3D shape generator $G_S$ to synthesize a 3D shape $S$, which we use as an explicit 3D shape prior. (b) In the second optimization stage, we use the 3D shape prior $S$ to initialize a neural radiance field and further optimize it with the CLIP guidance to synthesize 3D content consistent with the input text prompt $y$.

shown unprecedented performance recently. DALL-E [52] and CogView [11] use a autoregressive model [13, 52] architecture, while GLIDE [40] adopts a diffusion model [23, 64] architecture incorporated with CLIP [50] guidance or classifier-free guidance [24]. DALL-E-2 [51] leverages a diffusion prior network to translate CLIP text embeddings to CLIP image embeddings and an unCLIP decoder network to synthesize images from CLIP image embeddings. Imagen [58] and Stable Diffusion [55] utilize a large pre-trained language model to guide the sampling process of a diffusion model in pixel space and vector-quantized discrete latent space, respectively.

**Zero-Shot Text-to-3D.** Zero-shot [25,28,31,33,47,59] text-to-3D methods leverage the joint text-image modeling ability of pre-trained vision-language models, *e.g.*, CLIP [50], to eliminate the demand for paired text-3D data, which is hard to acquire. CLIP-Forge [59] trains a normalizing flow [12,54] model to translate CLIP image embeddings to VAE [29] shape embeddings, which replaces CLIP image embeddings with text embeddings at the inference phase to generate shapes. To close the gap between the text and shape modality, a fresh work named ISS [33] trains a mapper to map the CLIP image embedding to the shape latent feature of a pre-trained single-view reconstruction (SVR) network [42], then it finetunes the mapper by taking CLIP text embeddings as input. DreamFields [25] and CLIP-Mesh [28] are pioneering work that explores zero-shot 3D content creation with only CLIP guidance. The former optimizes a random-initialized NeRF while the latter optimizes a spherical template mesh as well as random texture and normal maps. PureCLIPNeRF [31] further improves DreamFields with grid-based representation [65] and more diverse image augmentations. Recently, DreamFusion [47] gains a lot of popularity in the research community with

its impressive results. Powered by a pre-trained 2D text-to-image model, it can generate high-fidelity relightable 3D objects with a probability density distillation loss.

## 3. Method

Given an input text prompt $y$, we aim to generate 3D content consistent with the description of $y$. As shown in Fig. 2, our text-guided 3D synthesis framework consists of two stages. In the first stage (Sec. E), we acquire an explicit 3D shape prior $S$ with a text-guided 3D shape generation process. The text-guided shape generation process is further split into a text-to-image phase using a fine-tuned Stable Diffusion model [55] $G_I$ and an image-to-shape phase using a shape embedding mapping network $G_M$ together with a high-quality 3D shape generator $G_S$. In the second stage, we then leverage the 3D prior $S$ to initialize a neural radiance field [38] and optimize it with the CLIP [50] guidance to obtain the final 3D content. Our framework only requires a collection of untextured 3D shapes $\{S\}_{i=1}^{N}$ without any text labels for training the 3D generator $G_S$, and the fine-tuning process of $G_I$ converges very fast.

### 3.1. CLIP-Guided 3D Optimization with explicit 3D Shape Prior

**Background: 3D Optimization with CLIP Guidance.** CLIP [50] is a powerful vision-language model consisting of a text encoder $E_T$ and an image encoder $E_I$. By maximizing the cosine similarity between the text embedding and the image embedding encoded by $E_T$ and $E_I$ respectively on a large-scale paired text-image training dataset, CLIP aligns the text and image modalities in a shared latent embedding space.

Previous work [25, 28, 31] leverages the power of

CLIP [50] to synthesize 3D content from text. Starting from a randomly-initialized 3D representation parameterized by $\theta$, they render images from multiple viewpoints and optimize $\theta$ by minimizing the CLIP similarity loss between the rendered image $\mathcal{R}(\boldsymbol{v}_i; \theta)$ and the text prompt $y$:

$$\mathcal{L}_{CLIP} = -E_I(\mathcal{R}(\theta; \boldsymbol{v}_i)^T E_T(y), \qquad (1)$$

where $\mathcal{R}$ denotes the rendering process and $\boldsymbol{v}_i$ denotes the rendering viewpoint at the $i$-th optimization step. Specifically, DreamFields [25] and PureCLIPNeRF [31] adopt neural radiance fields [38] (NeRF) as the 3D representation $\theta$, while CLIP-Mesh [28] adopts a spherical template mesh associated with texture and normal maps.

**Observations and Motivations.** Despite these CLIP-guided optimization methods can produce impressive results, we observe that they usually fail to generate accurate and detailed 3D structures that conform to the corresponding text description. As Fig. 1 illustrates, we utilize those methods to synthesize 3D contents with text prompts that contain common objects. The results show distortion artifacts and look a bit strange, which has a negative impact on the visual quality and restrains their deployment in real-world applications.

We consider there are two main reasons for the failure of previous work in generating accurate and faithful objects: **(i)** The optimization process starts from a *randomly-initialized 3D representation* without any explicit 3D shape prior. Thus, the models need to imagine the scene from scratch, which is very challenging. **(ii)** The CLIP loss in Eq. (1) concentrates more on the global consistency between the rendered image and the text prompt, and cannot provide strong and precise supervision on the synthesized 3D structure. Therefore, the structure of the optimized 3D representation is highly unconstrained.

**Optimization with 3D Shape Prior as Initialization.** To overcome the aforementioned problem and synthesize more faithful 3D content, we propose to adopt a text-guided shape synthesis process to produce a *high-quality 3D shape $S$* from the input text prompt $y$, and then leverage it as an explicit "3D shape prior" to initialize the CLIP-guided 3D optimization process. As Fig. 2 shows, for the input text prompt *"a park bench overgrown with vines"*, we first synthesize *"a park bench"* 3D shape without textures in the text-to-shape stage, and then utilize it as the initialization of neural radiance fields and then optimize with the full prompt following previous works.

Our optimization process adopts an efficient NeRF representation [38], *i.e.*, DVGO [65], which represents NeRF with a density voxel grid $\boldsymbol{V}_{density} \in \mathbb{R}^{N_x \times N_y \times N_z}$ and a shallow MLP network $f_{rgb}$ for color prediction. Given a 3D shape $S$ produced by the text-to-shape process represented as an SDF grid $\tilde{\boldsymbol{V}}_{sdf} \in \mathbb{R}^{N_x \times N_y \times N_z}$, which will be introduced in Sec. E, we first transform it into the density voxel grid $\boldsymbol{V}_{density}$ using the following equations [44, 65, 76]:

$$\boldsymbol{\Sigma} = \frac{1}{\beta} \operatorname{sigmoid}\left(-\frac{\tilde{\boldsymbol{V}}_{sdf}}{\beta}\right), \qquad (2)$$

$$\boldsymbol{V}_{density} = \max(0, \operatorname{softplus}^{-1}(\boldsymbol{\Sigma})) \qquad (3)$$

where $\operatorname{sigmoid}(x) = 1/(1 + e^{-x})$ and $\operatorname{softplus}^{-1}(x) = \log(e^x - 1)$. Eq. (9a) converts SDF values to density for volume rendering, where $\beta > 0$ is a hyper-parameter controlling the sharpness of the shape boundary, and smaller $\beta$ leads to a sharper shape boundary. We set $\beta = 0.05$ in our experiments. Eq. (9b) transform the density into pre-activated density. We follow the practice of DVGO and adopt this representation for the density voxel grid. We clamp the minimum density value of outside the shape prior as 0 to ensure the distribution of the accumulated transmittance is the same as DVGO.

With the density grid $\boldsymbol{V}_{density}$ initialized by the explicit 3D shape $S$ and the color MLP $f_{rgb}$ initialized randomly, we render images $\mathcal{R}(\boldsymbol{V}_{density}, f_{rgb}; \boldsymbol{v}_i)$ from viewpoint $\boldsymbol{v}_i$ and optimize $\boldsymbol{V}_{density}$ and $f_{rgb}$ with the CLIP loss in Eq. (1). Following DreamFields [25] and PureCLIPN-eRF [31], we perform background augmentations for the rendered images and leverage the transmittance loss introduced by [25] to reduce noise and spurious density. Besides, since CLIP loss cannot provide accurate geometrical supervision, the 3D shape prior may be disturbed and forgotten gradually, thus we also adopt a shape-prior-preserving loss to preserve the global structure of the 3D shape prior:

$$\mathcal{L}_{prior} = -\sum_{x,y,z} \mathbb{1}(\tilde{\boldsymbol{V}}_{sdf} < 0) \cdot \operatorname{alpha}(\boldsymbol{V}_{density}), \qquad (4)$$

where $\mathbb{1}(\cdot)$ is the indicator function and $\operatorname{alpha}(\cdot)$ transforms the density into the opacity representing the probability of termination at each position in volume rendering.

By initializing the NeRF representation with the explicit 3D shape prior, we give extra knowledge on how the 3D content should look like and prevent the model from imagining from scratch and generating "adversarial contents" that have high CLIP scores but low visual quality. Based on this initialization, the CLIP-guided optimization process further provides flexibility and is able to synthesize *more diverse structures and textures*.

### 3.2. Stable-Diffusion-Assisted Text-to-Shape Generation as 3D Shape Prior

To acquire 3D shape priors, we require a text-guided shape generation scheme. Text-to-shape is inherently a challenging task due to the lack of large paired text-shape dataset. Previous work [33, 59] typically trains an image-to-shape module using *rendered* images first and then bridges the text and image modalities using the CLIP embedding space.

Specifically, CLIP-Forge [59] trains a volumetric shape auto-encoder and a normalizing flow network to map CLIP image embeddings of ShapeNet renderings to shape embeddings. At test time, it replaces the image embeddings with CLIP text embeddings. However, the shape auto-encoder has difficulty generating high-quality and diversified 3D shapes, and directly feeding CLIP text embeddings to the shape generation network trained on CLIP image embeddings suffers from the gap between the CLIP text and image embeddings. ISS [33] trains a mapper network to map the CLIP image embeddings of shape renderings to the latent space of a pre-trained single-view reconstruction (SVR) model, and fine-tunes the mapper at test time by maximizing the CLIP similarity between the input text prompt and the images rendered from synthesized shapes. Although the test-time fine-tuning strategy alleviates the gap between the CLIP text and image embeddings, it is tedious to fine-tune the mapper for each text prompt.

Different from the aforementioned approaches which connect the text and image modalities in CLIP embedding space, we use a powerful text-to-image diffusion model to bridge the two modalities directly. Specifically, we adopt the text-to-image diffusion models to synthesize an image from the input text prompt, and then feed the image into an image-to-shape module to generate *high-quality* 3D shapes. This pipeline is more concise and naturally eliminates the gap between text and image CLIP embeddings. However, it brings a new domain gap between the images synthesized by the text-to-image model and the shape renderings used to train the image-to-shape module. We will introduce a novel technique to alleviate this gap in Sec. 3.3.

**Text-to-Image Diffusion Model.** Diffusion models [23,64] are generative models that are trained to invert a diffusion process. Given a sample from the data distribution $x_0 \sim q(x_0)$, the forward diffusion process gradually adds Gaussian noise to it in $T$ timesteps: $x_t = \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \epsilon_{t-1}, t = 1, 2, \ldots, T$, where $\alpha_t$ defines the noise level and $\epsilon_{t-1}$ denotes the noise added at timestep $t - 1$. A neural network $\epsilon_\theta$ is trained to reverse the diffusion process by estimating the added noise at each timestep. At inference time, we can sample from the diffusion model iteratively through a sampler. Text-to-image diffusion models are conditional diffusion models. Given a text prompt $y$ and let $c_\theta$ denote a text encoder that maps $y$ into a conditioning vector, the training objective of text-to-image diffusion models is written as:

$$\mathcal{L}_{diffusion} = \mathbb{E}_{x_0, t, \epsilon_t, y} \|\epsilon_t - \epsilon_\theta(x_t, t, c_\theta(y))\|_2^2 \quad (5)$$

In this work, we use an open source text-to-image diffusion model, *i.e.*, Stable Diffusion[†] [55]. Stable Diffusion utilizes a CLIP ViT-L/14 text encoder as $c_\theta$ and is trained on the

---

[†] https://github.com/CompVis/StableDiffusion

large-scale LAION-5B dataset [60]. It has superior ability in generating diverse and imaginative images in various styles from heterogeneous text prompts.

**High-quality 3D generator.** High-quality 3D shapes are highly desired and more precise 3D geometric shapes can provide a better initialization for optimization. CLIP-Forge [59] utilizes a volumetric shape auto-encoder, which has difficulty in generating plausible 3D shapes. To provide high-quality 3D priors, we adopt a state-of-the-art 3D generative model architecture, *i.e.*, SDF-StyleGAN [79]. SDF-StyleGAN is a StyleGAN2-like generative architecture that maps a random noise $z \sim \mathcal{N}(0, I)$ to a latent shape embedding $e_S \in \mathcal{W}$ and synthesizes a 3D feature volume $F_V$ - an implicit representation of the generated shape. We can query the SDF value at arbitrary position $x$ by feeding the interpolated feature from $F_V$ at $x$ into a jointly trained MLP network. Different from the original SDF-StyleGAN that trains one network for one shape category, we improve it to train *one* 3D shape generator $G_S$ on 13 categories of the ShapeNet [4] dataset, so that the text-to-shape process has more flexibility.

**Shape Embedding Mapping Network.** To bridge the image and shape modalities, we further train a shape embedding mapping network $G_M$. We first use $G_S$ to randomly generate a large set of 3D shapes $\{S^i\}_{i=1}^N$ and their corresponding shape embeddings $\{e_S^i\}_{i=1}^N$. Then we render each shape from $K$ viewpoints to obtain shape renderings $\{I_r^j\}_{j=1}^M, M = NK$. After transforming the shape renderings to image embeddings $\{e_I^j\}_{j=1}^M$ with the CLIP image encoder $E_I$, we obtain a paired image-shape embedding dataset $\{(e_I^j, e_S^j)\}_{j=1}^M$. Finally, we use this dataset to train a conditional diffusion model $G_M$ which can synthesize shape embeddings from image embeddings of shape renderings.

To prepare the dataset for training $G_M$, we generate $N = 64000$ shapes using the 3D generator $G_S$ and render $K = 24$ views for each shape with azimuth angles in the range of $[-90°, 90°]$ and elevation angles in the range of $[20°, 30°]$. We adopt an SDF renderer [26] since we represent the synthesized shapes with SDF grids.

### 3.3. Fine-tuning Stable Diffusion for Rendering-Style Image Generation

As we have stated in Sec. E, we use Stable Diffusion to directly bridge the text and image modalities for text-to-shape generation. However, the image-to-shape module is trained on shape renderings, which have a large style domain gap from the images synthesized by Stable Diffusion. Previous works [33, 59] has tried simply combining a text-to-image model [40, 52, 81] and an image-to-shape model for text-to-shape generation. However, this approach suffers from the aforementioned style domain gap, leading to distorted geometric structures and degraded performance.
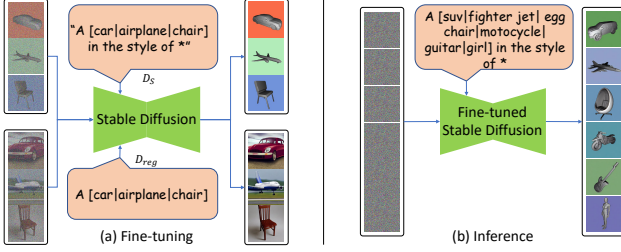
Figure 3. Fine-tuning Stable Diffusion into a stylized generator. (a): We use a dataset $D_S$ containing ShapeNet renderings and text descriptions in the format of *"a CLS in the style of *"* to fine-tune the text embedding $\boldsymbol{v}_*$ of the placeholder token $*$ and the weights of the diffusion model, another dataset $D_{reg}$ which contains images synthesized by the original Stable Diffusion with text prompts in the format of *"a CLS"* is also used for regularization. (b): with the fine-tuned Stable Diffusion, we can synthesize images that meets the style of ShapeNet renderings by adding a postfix *"in the style of *"* to the text prompt.

Inspired by recent work on controllable text-to-image synthesis, *e.g.*, textual inversion [14] and DreamBooth [57], we propose an approach to address this domain gap by fine-tuning Stable Diffusion into a stylized generator. Our core idea is to give the original Stable Diffusion the ability to synthesize images that meets the style of shape renderings we used to train the image-to-shape module as in Sec. E. Then we can seamlessly feed the generated *stylized* images to the image-to-shape module without suffering from the domain gap when using images synthesized by the original Stable Diffusion.

**Fine-tuning Process.** The fine-tuning process is shown in Fig. 8. To fine-tune Stable Diffusion, we need a dataset of shape renderings and associated *stylized* text prompts. For each shape $S$ in the ShapeNet dataset, we render a set of shape renderings $\{\boldsymbol{I}_S^j\}_{j=1}^{N_S}$. Then for each rendering $\boldsymbol{I}_S^j$, we associate it with a *stylized* text prompt $y_S^j$ in the format of *"a CLS in the style of *"*, where *CLS* is replaced by the shape category name and $*$ is a placeholder token whose text embedding needs to be optimized. For example, for an image rendered from a chair shape, the associated text prompt is *"a chair in the style of *"*. We then use the paired dataset $D_S = \{(\boldsymbol{I}_S^j, y_S^j)\}_{j=1}^{N_S}$ to fine-tune Stable Diffusion by minimizing the objective in Eq. (5).

During fine-tuning, we freeze the CLIP text encoder of Stable Diffusion, and optimize two objectives: (i) the text embedding of the placeholder token $*$, denoted as $\boldsymbol{v}_*$, and (ii) the parameters $\theta$ of the diffusion model $\epsilon_\theta$. By optimizing the text embedding $\boldsymbol{v}_*$, we aim to learn a virtual word that is not in the vocabulary of the text encoder but can best describe the style of the rendered images. And fine-tuning the parameters $\theta$ of the diffusion model can further help to capture the style more precisely since it is hard to control the synthesis process of Stable Diffusion at only the language level. In our experiments, the fine-tuning pro-

cess converges stably in around 2000 optimization steps, which only takes about 40 minutes on a single Tesla A100 GPU. We show some synthesized results using the fine-tuned model in Fig. 8.

**Dataset Scale and Background Augmentation.** We empirically find two techniques essential to make the fine-tuned model stably synthesize stylized images as we wish. Firstly, unlike textual inversion [14] or DreamBooth [57] utilizing only $3-5$ images, fine-tuning with a larger set of shape renderings containing thousands of images helps the model capture the style more precisely. Secondly, fine-tuning Stable Diffusion using shape renderings with a pure-white background would result in a chaotic and uncontrollable background during inference, while augmenting the shape renderings with random solid-color backgrounds can make the fine-tuned model synthesize images with solid-color backgrounds stably, thus we can remove the background easily if needed. Please refer to the supplementary material for more details.

## 4. Experiments

In this section, we evaluate the effectiveness of our proposed text-guided 3D synthesis framework. We first compare our text-guided 3D synthesis and results with state-of-the-art methods (Sec. 4.1). Then we demonstrate the effectiveness of our proposed Stable-Diffusion-assisted text-to-shape generation approach (Sec. 4.2). We further conduct ablation studies to evaluate the effectiveness of several crucial parts of our framework (Sec. 4.3).

**Dataset.** Our framework only requires a set of untextured 3D shapes to train the 3D shape generator $G_S$. Specifically, we use 13 categories of ShapeNet [4] and follow the data preprocessing procedure of Zheng *et al.* [79] to produce $128^3$ SDF grids from the raw meshes for training the 3D generator. When fine-tuning Stable Diffusion, we render the SDF grids with the SDF renderer [26] to obtain a shape rendering dataset.

**Implementation details.** Our 3D generator $G_S$ adopts the architecture of SDF-StyleGAN [79], and the diffusion-model-based shape embedding mapping network $G_M$ is based on an open-source DALLE-2 implementation$^\ddagger$. To train the shape embedding mapping network $G_M$, we use the CLIP ViT-B/32 image encoder to extract image embedding from shape renderings. Our stylized text-to-image generator $G_I$ is fine-tuned from the *"sd-v1-4-full-ema"* checkpoint of Stable Diffusion. In the optimization stage. the learning rates for the density grid $\boldsymbol{V}_{density}$ and color prediction network $f_{rgb}$ are set to $5 \times 10^{-1}$ and $5 \times 10^{-3}$ respectively, and the CLIP ViT-B/16 encoder is adopted as the guidance model. For each text prompt, we optimize for 5000 steps, while previous NeRF-based 3D synthesis meth-

---

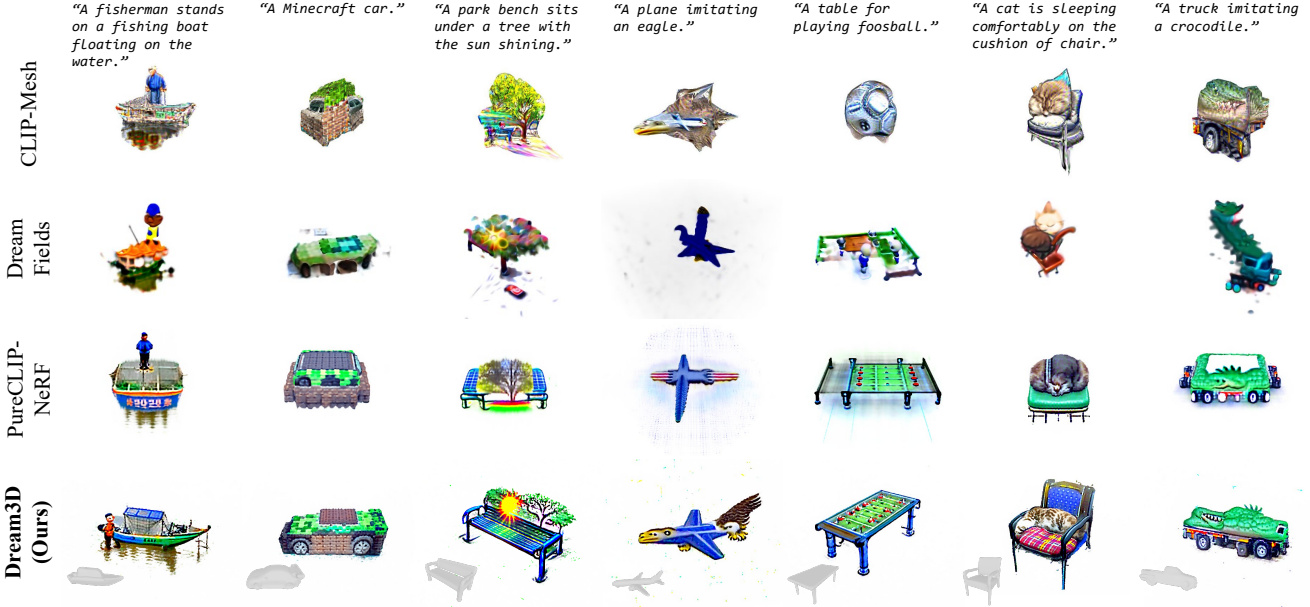$^\ddagger$https://github.com/lucidrains/DALLE2-pytorch

Figure 4. Qualitative comparison on text-guided 3D synthesis. For each of our results (the last row), we also visualize the 3D shape prior used to initialize the CLIP-guided optimization process below it.

ods [25, 31] typically require 10000 steps or more.

**Evaluation Metrics.** For the main results of our framework, *i.e.*, text-guided 3D content synthesis, we report the CLIP retrieval precision on a hand-crafted diversified text prompt dataset containing objects, please refer to the supplementary material for the details of the dataset. This metric measures the fraction of generated images that the CLIP encoder associates with the text prompt used to generate it. For the first-stage text-to-shape generation, we employ Fréchet Inception Distance (FID) [22] to measure the shape generation quality.

### 4.1. Text-Guided 3D Synthesis

We compare our method with three state-of-the-art baseline methods on the task of text-guided 3D synthesis, *i.e.*, DreamFields [25], CLIP-Mesh [28], and PureCLIPNeRF [31]. We test with the official implementations and the default configurations for all baseline methods. Specifically, we use the medium-quality configuration of DreamFields and the implicit architecture variant of PureCLIPNeRF on account of its better performance. And the default optimization steps for DreamFields, CLIP-Mesh, and Pure-CLIPNeRF are 10000, 5000, and 15000 respectively.

**Quantitative Results.** We report the CLIP retrieval precision metrics in Tab. 1. To be noted, all the baseline methods and ours leverage the CLIP ViT-B/16 encoder for optimization, and both the CLIP ViT-B/16 and CLIP ViT-B/32 encoder are used as the retrieval model. As Tab. 1 shows, our method achieves the highest CLIP R-Precision with both retrieval models. Besides, we can also observe that our

framework shows a significantly smaller performance gap between the two retrieval models than baseline methods. Thanks to the 3D shape prior, our method starts the optimization from a better point, thus alleviating the adversarial generation problem that focuses on obtaining high score from the guidance CLIP model while ignoring the visual quality. Therefore, our method shows more robust performance between different CLIP models.

**Qualitative Results.** We show the qualitative comparison in Fig. 12. It is easy to observe that the baseline methods [25, 28, 31] have difficulties in generating precise and plausible 3D objects, showing obvious distortions and unrealistic visuals. Specifically, the synthesis results of Dream-Fields are often blurry and diffuse, while PureCLIPNeRF tends to synthesize symmetric objects. CLIP-Mesh has difficulty in generating complicated visual effects due to its explicit mesh representation. In comparison, our method generates higher-quality 3D structures thanks to the injection of explicit 3D shape priors.

### 4.2. Text-to-Shape Generation

The research on zero-shot text-to-shape generation is still limited, thus we only compare our text-to-shape generation approach with CLIP-Forge [59]. We measure the shape generation quality with the Fréchet Inception Distance (FID). Concretely, we first synthesize 3 shapes for each prompt from a dataset consisting of 233 text prompts provided by CLIP-Forge, and then render 5 images for each synthesized shape and compare them with a set of ground truth ShapeNet renderings to compute the FID. As Tab. 5
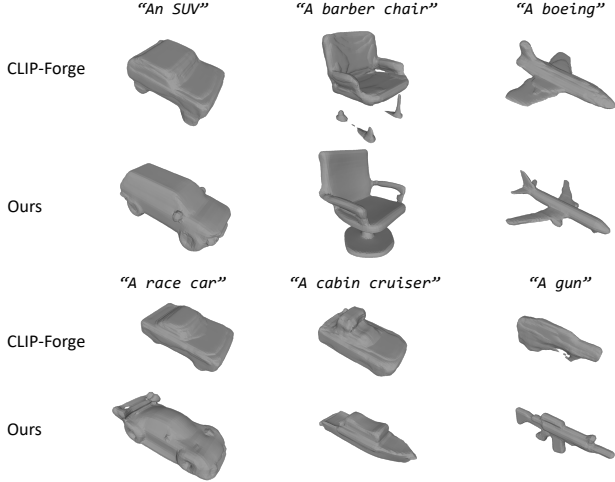
"An SUV"    "A barber chair"    "A boeing"

CLIP-Forge

Ours

"A race car"    "A cabin cruiser"    "A gun"

CLIP-Forge

Ours

Figure 5. Text-guided 3D shape generation results. All the visualized meshes are extracted at the resolution of $64^3$. We can observe that our method can generate significantly more plausible 3D shapes benefitting from the high-quality 3D shape generator.

| Method | CLIP R-Precision ↑ | |
| --- | --- | --- |
| | ViT-B/32 | ViT-B/16 |
| DreamFields [25] | 63.24 | 92.65 |
| CLIP-Mesh [28] | 75.00 | 91.18 |
| PureCLIPNeRF [31] | 73.53 | 88.24 |
| Ours w/o 3D prior | 75.47 | 94.34 |
| Ours | **85.29** | **98.53** |

Table 1. Quantitative comparison on Text-guided 3D content synthesis. All the methods utilize CLIP ViT/16 as the guidance model during optimization, while we use two different CLIP models to calculate the CLIP retrieval precision.

| Method | FID ↓ |
| --- | --- |
| CLIP-Forge [59] | 112.38 |
| Ours w/o text-to-image | 58.36 |
| Ours w/o fine-tuning SD | 61.88 |
| Ours | **40.83** |

Table 2. Quantitative comparison on text-to-shape generation.

shows, we achieve a lower FID than CLIP-Forge. CLIP-Forge leverages a volumetric shape auto-encoder to generate 3D shapes, which suffers from poop shape generation capability and is hard to generate plausible 3D shapes, the qualitative results in Fig. 10 shows. A high-quality 3D shape prior is also beneficial to the optimization process as a good initialization.

## 4.3. Ablation Studies

**Effectiveness of Fine-tuning Stable Diffusion.** Our framework leverages a fine-tuned Stable Diffusion to bridge the text and image modalities. To validate its effectiveness, we use the original Stable Diffusion model to generate images from the text prompts CLIP-Forge [59], and then use these images to generate 3D shapes with the image-to-shape module. We can observe an FID performance drop in the third
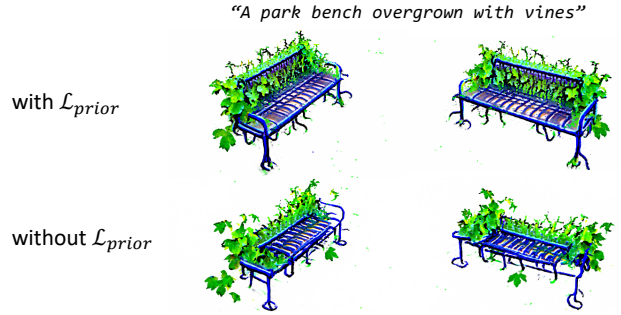


"A park bench overgrown with vines"

with $\mathcal{L}_{prior}$

without $\mathcal{L}_{prior}$

Figure 6. The effect of the 3D prior preserving loss $\mathcal{L}_{prior}$.

row of Tab. 5 indicating that this strategy would degrade the shape generation performance. Besides, we also test directly using text embedding to generate shape embeddings with $G_M$, which also leads to a performance drop as the second row of Tab. 5 shows.

**Effectiveness of 3D Shape Prior.** To validate the effectiveness of the 3D shape prior, we remove the first stage from our framework and optimize from scratch with the same text prompts as in Sec. 4.1, and evaluate the CLIP retrieval precision. As Tab. 1 shows, optimizing without 3D shape prior degrades the performance significantly, which demonstrates its effectiveness.

**Effectiveness of Shape Prior Preserving Loss.** The 3D prior preserving loss $\mathcal{L}_{prior}$ in Eq. (4) aims to strengthen the 3D prior in the optimization process in case the prior is gradually disturbed and discarded. Using the text prompt *"A park bench overgrown with vines"*, we first synthesize a park bench shape as 3D prior, then we perform optimization with and without the prior preserving loss respectively for comparison, and visualize the optimization results in Fig. 6. We can observe that using $\mathcal{L}_{prior}$ during optimization helps maintain the structure of the 3D prior shape, while discarding $\mathcal{L}_{prior}$ would disturb the initial shape and lead to distortion and discontinuity artifacts.

## 5. Limitations and Future Work

Our framework depends on a fine-tuned Stable Diffusion model to synthesize stylized images to feed into the image-to-shape module. Despite the strong generation capability of Stable Diffusion, we cannot constrain it to avoid generating shape images that are out of the distribution of the 3D shape generator, since Stable Diffusion is trained on a mega-scale text-image dataset while the 3D shape generator can only generate a limited amount of shapes. Besides, the quality of text-to-shape synthesis in our framework highly depends on the generation capability of the 3D generator. In future work, we hope to introduce stronger 3D priors into our framework to make it work on more extensive object categories.

## 6. Conclusion

In this paper, we present Dream3D, a text-guided 3D synthesis framework to generate diversified and imaginative 3D content from texts. Our framework leverages text-to-shape synthesis to inject explicit 3D priors into the CLIP-guided optimization process for more plausible 3D structures. For the text-to-shape generation, we present a simple yet effective approach that directly use a fine-tuned text-to-image diffusion model to connect the text and image modalities. Extensive experiments demonstrate the capability of our method of generating imaginative 3D content with better visual quality and shape accuracy than previous work.

## References

[1] Ruojin Cai, Guandao Yang, Hadar Averbuch-Elor, Zekun Hao, Serge Belongie, Noah Snavely, and Bharath Hariharan. Learning gradient fields for shape generation. In *European Conference on Computer Vision*, pages 364–381. Springer, 2020. 2

[2] 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, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16123–16133, June 2022. 2

[3] Eric R. Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. 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 (CVPR)*, pages 5799–5809, June 2021. 2

[4] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 5, 6, 13

[5] Kevin Chen, Christopher B Choy, Manolis Savva, Angel X Chang, Thomas Funkhouser, and Silvio Savarese. Text2shape: Generating shapes from natural language by learning joint embeddings. In *Asian conference on computer vision*, pages 100–116. Springer, 2018. 2

[6] Wenzheng Chen, Huan Ling, Jun Gao, Edward Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. *Advances in Neural Information Processing Systems*, 32, 2019. 2

[7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2

[8] Zezhou Cheng, Menglei Chai, Jian Ren, Hsin-Ying Lee, Kyle Olszewski, Zeng Huang, Subhransu Maji, and Sergey Tulyakov. Cross-modal 3d shape generation and manipula-tion. In *European Conference on Computer Vision (ECCV)*, 2022. 2

[9] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. 19

[10] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10673–10683, June 2022. 2

[11] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835, 2021. 1, 3

[12] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *International Conference on Learning Representations*, 2017. 3

[13] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12873–12883, June 2021. 3

[14] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. 2, 6

[15] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from images. In *Advances In Neural Information Processing Systems*, 2022. 2

[16] Lin Gao, Tong Wu, Yu-Jie Yuan, Ming-Xian Lin, Yu-Kun Lai, and Hao Zhang. Tm-net: Deep generative networks for textured meshes. *ACM Transactions on Graphics (TOG)*, 40(6):263:1–263:15, 2021. 2

[17] Lin Gao, Jie Yang, Tong Wu, Yu-Jie Yuan, Hongbo Fu, Yu-Kun Lai, and Hao Zhang. Sdm-net: Deep generative network for structured deformable mesh. *ACM Transactions on Graphics (TOG)*, 38(6):1–15, 2019. 2

[18] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. 2

[19] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d aware generator for high-resolution image synthesis. In *International Conference on Learning Representations*, 2022. 2

[20] Kunal Gupta and Manmohan Chandraker. Neural mesh flow: 3d manifold mesh generation via diffeomorphic flows. In *Advances in Neural Information Processing Systems*, volume 33, pages 1747–1758, 2020. 2

[21] Paul Henderson, Vagia Tsiminaki, and Christoph H. Lampert. Leveraging 2d data to learn textured 3d mesh genera-

tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2

[22] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 7, 15

[23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 3, 5

[24] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 3

[25] 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 (CVPR)*, pages 867–876, June 2022. 1, 2, 3, 4, 7, 8, 15

[26] Yue Jiang, Dantong Ji, Zhizhong Han, and Matthias Zwicker. Sdfdiff: Differentiable rendering of signed distance fields for 3d shape optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2, 5, 6

[27] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2

[28] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Popa Tiberiu. Clip-mesh: Generating textured meshes from text using pretrained image-text models. December 2022. 1, 3, 4, 7, 8, 15

[29] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3

[30] Marian Kleineberg, Matthias Fey, and Frank Weichert. Adversarial generation of continuous implicit shape representations. *arXiv preprint arXiv:2002.00349*, 2020. 2

[31] Han-Hung Lee and Angel X Chang. Understanding pure clip guidance for voxel grid nerf models. *arXiv preprint arXiv:2209.15172*, 2022. 1, 2, 3, 4, 7, 8, 15

[32] Jun Li, Kai Xu, Siddhartha Chaudhuri, Ersin Yumer, Hao Zhang, and Leonidas Guibas. Grass: Generative recursive autoencoders for shape structures. *ACM Transactions on Graphics (TOG)*, 36(4):1–14, 2017. 2

[33] Zhengzhe Liu, Peng Dai, Ruihui Li, Xiaojuan Qi, and Chi-Wing Fu. Iss: Image as stetting stone for text-guided 3d shape generation. *arXiv preprint arXiv:2209.04145*, 2022. 1, 2, 3, 4, 5, 19

[34] Zhengzhe Liu, Yi Wang, Xiaojuan Qi, and Chi-Wing Fu. Towards implicit text-guided 3d shape generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17896–17906, June 2022. 2

[35] Andrew Luo, Tianqin Li, Wen-Hao Zhang, and Tai Sing Lee. Surfgen: Adversarial 3d shape synthesis with explicit surface discriminators. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16238–16248, October 2021. 2

[36] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2837–2845, June 2021. 2

[37] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 19

[38] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2, 3, 4, 14

[39] Charlie Nash, Yaroslav Ganin, S. M. Ali Eslami, and Peter W. Battaglia. Polygen: An autoregressive generative model of 3d meshes. *ICML*, 2020. 2

[40] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 1, 3, 5

[41] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11453–11464, June 2021. 2

[42] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2, 3, 19

[43] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. Stylesdf: High-resolution 3d-consistent image and geometry generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13503–13513, June 2022. 2

[44] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. Stylesdf: High-resolution 3d-consistent image and geometry generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13503–13513, 2022. 4, 15

[45] Dario Pavllo, Jonas Kohler, Thomas Hofmann, and Aurelien Lucchi. Learning generative models of textured 3d meshes from real-world images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13879–13889, October 2021. 2

[46] Dario Pavllo, Graham Spinks, Thomas Hofmann, Marie-Francine Moens, and Aurelien Lucchi. Convolutional generation of textured 3d meshes. *Advances in Neural Information Processing Systems*, 33:870–882, 2020. 2

[47] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv*, 2022. 1, 3

[48] Albert Pumarola, Stefan Popov, Francesc Moreno-Noguer, and Vittorio Ferrari. C-flow: Conditional generative flow

models for images and 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2

[49] Tingting Qiao, Jing Zhang, Duanqing Xu, and Dacheng Tao. Mirrorgan: Learning text-to-image generation by redescription. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2

[50] 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. 3, 4

[51] 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, 13

[52] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 1, 3, 5

[53] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *International conference on machine learning*, pages 1060–1069. PMLR, 2016. 2

[54] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015. 3

[55] 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 (CVPR)*, pages 10684–10695, June 2022. 1, 2, 3, 5

[56] Shulan Ruan, Yong Zhang, Kun Zhang, Yanbo Fan, Fan Tang, Qi Liu, and Enhong Chen. Dae-gan: Dynamic aspect-aware gan for text-to-image synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13960–13969, October 2021. 2

[57] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arXiv:2208.12242*, 2022. 2, 6

[58] 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

[59] Aditya Sanghi, Hang Chu, Joseph G. Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshan. Clip-forge: Towards zero-shot text-to-shape generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18603–18613, June 2022. 1, 2, 3, 4, 5, 7, 8, 15

[60] 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. 2, 5

[61] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 2

[62] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020. 2

[63] Dong Wook Shu, Sung Woo Park, and Junseok Kwon. 3d point cloud generative adversarial network based on tree structured graph convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 15

[64] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 3, 5

[65] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5459–5469, June 2022. 3, 4, 14, 15

[66] Hongchen Tan, Xiuping Liu, Xin Li, Yi Zhang, and Baocai Yin. Semantics-enhanced adversarial nets for text-to-image synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 2

[67] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 2

[68] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019. 16

[69] Zihao Wang, Wei Liu, Qian He, Xinglong Wu, and Zili Yi. Clip-gen: Language-free training of a text-to-image generator with clip. *arXiv preprint arXiv:2203.00386*, 2022. 1

[70] Rundi Wu and Changxi Zheng. Learning to generate 3d shapes from a single example. *ACM Transactions on Graphics (TOG)*, 41(6), 2022. 2

[71] Rundi Wu, Yixin Zhuang, Kai Xu, Hao Zhang, and Baoquan Chen. Pq-net: A generative part seq2seq network for 3d shapes. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2

[72] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2

[73] Xudong Xu, Xingang Pan, Dahua Lin, and Bo Dai. Generative occupancy fields for 3d surface-aware image synthe-

sis. *Advances in Neural Information Processing Systems*, 34:20683–20695, 2021. 2

[74] Yang Xue, Yuheng Li, Krishna Kumar Singh, and Yong Jae Lee. Giraffe hd: A high-resolution 3d-aware generative model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18440–18449, June 2022. 2

[75] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 2

[76] 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. 4, 15

[77] Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten Kreis. Lion: Latent point diffusion models for 3d shape generation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2

[78] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 2

[79] Xin-Yang Zheng, Yang Liu, Peng-Shuai Wang, and Xin Tong. Sdf-stylegan: Implicit sdf-based stylegan for 3d shape generation. In *Comput. Graph. Forum (SGP)*, 2022. 2, 5, 6, 13, 14

[80] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5826–5835, October 2021. 2

[81] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. Towards language-free training for text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17907–17917, June 2022. 5

## A. Details of 3D Generator $G_S$

We adopt the architecture of SDF-StyleGAN [79] as our 3D generator. As Fig. 7 shows, it maps a random noise $z \sim \mathcal{N}(0, I)$ to a latent shape embedding $e_S \in \mathcal{W}$ and synthesizes a 3D feature volume $F_V$, which is an implicit representation of the generated shape. We can query the SDF value at arbitrary position $x$ by feeding the interpolated feature from $F_V$ at $x$ into a jointly trained MLP network. During training, a global discriminator and a local discriminator are used simultaneously to supervise the generated SDF grids at the coarse and fine level respectively. Different from the original SDF-StyleGAN that trains one network for *one shape category*, we train *one* 3D shape generator $G_S$ on *13 categories* of the ShapeNet [4] dataset to enlarge the shape generation capability.

## B. Details of Shape Embedding Mapping Network $G_M$

The shape embedding mapping network $G_M$ is a diffusion-model-based generative network that can generate shape embeddings $e_S$ from the CLIP image embeddings $e_I$ of shape renderings. The network architecture and training strategy of $G_M$ are based on an open-source DALL-E-2 [51] implementation[§]. Specifically, $G_M$ is equivalent to the *diffusion prior network* in DALL-E-2 which generates CLIP image embeddings from CLIP text embeddings. Here we replace the input with CLIP image embeddings of shape renderings and the output with shape embeddings. We use the *DiffusionPrior* class in the codebase to implement $G_M$ and the *train_diffusion_prior.py* script to train $G_M$. The model and training hyperparameters are listed in Tab. 3.

## C. Details of Fine-tuning Stable Diffusion

In our framework, we connect the text and image modalities by fine-tuning the Stable Diffusion into a stylized generator with a set of shape renderings $\{I_S^j\}_{j=1}^{N_S}$ and give it the ability to synthesize images in the "rendering" style. In experiments, we find it crucial to utilize a large set of shape renderings for fine-tuning and to augment the backgrounds of the shape renderings with random colors.

We tried fine-tuning Stable Diffusion using shape renderings with three different types of backgrounds: 1) solid white background, 2) solid green background and 3) random-color background. We visualize the shape renderings used for fine-tuning and the images synthesized by the fine-tuned Stable Diffusion in Fig. 8. As Fig. 8a and Fig. 8b show, although shape renderings with solid-white or solid-green backgrounds can make the fine-tuned Stable Diffusion capture the "rendering" style of the object successfully, the backgrounds in the synthesized images are

| Model Parameter | Value | Training Parameter | Value |
|---|---|---|---|
| timesteps | 100 | iterations | 500,000 |
| beta_schedule | cosine | max_grad_norm | 0.5 |
| predict_x_start | True | batch_size | 1024 |
| cond_drop_prob | 0 | learning_rate | $1.1 \times 10^{-4}$ |
| dim | 512 | weight_decay | $6.02 \times 10^{-2}$ |
| depth | 6 | ema_beta | 0.9999 |
| dim_head | 64 | ema_update_every | 10 |
| heads | 8 | Adam $\beta_1, \beta_2$ | 0.9, 0.999 |

Table 3. Model details and training hyper-parameters of the shape embedding mapping network $G_M$.

| Background | FID $\downarrow$ |
|---|---|
| Solid white | 60.61 |
| Solid green | 71.04 |
| Random-color | **33.71** |

Table 4. The Fréchet Inception Distance (FID) between the shape renderings used for fine-tuning and the images synthesized by the fine-tuned Stable Diffusion.

out of control, *i.e.*, the fine-tuned Stable Diffusion fails to synthesize images with solid-color backgrounds. This will increase the difficulty of separating the foreground objects from the backgrounds and affect the stability of the subsequent image-to-shape generation since the shape embedding mapping network $G_M$ is trained on shape renderings with solid-color backgrounds. In comparison, augmenting the backgrounds of the shape renderings with random colors leads to a stable stylized generator that can synthesize solid-color-background images consistently, as Fig. 8c shows.

During fine-tuning, we indeed expect the Stable Diffusion model to capture two types of styles: 1) the "rendering" style of the foreground object and 2) the "solid-color" style of the background. Similar to the observation that the foreground "rendering" style requires a large set of rendered images to learn, we consider that a single-color background is too few to be recognized as a "solid-color background style" by the Stable Diffusion model, while showing a lot of different solid-color examples to the model can make it notice the solid-color background style and capture it during fine-tuning.

To better demonstrate the importance of the random-color background augmentation, we also evaluate the Fréchet Inception Distance (FID) between the shape renderings used for fine-tuning and the images synthesized by the fine-tuned Stable Diffusion in Tab. 4. For each type of background, we render 1000 images with that background for each ShapeNet [4] category, forming a shape rendering dataset containing 13000 images in total (denoting as $D_S$). Then we leverage $D_S$ to fine-tune the Stable Diffusion model for 5000 steps, and utilize the fine-tuned Stable Diffusion to synthesize 100 images for each shape category using the text prompt *"a CLS in the style of *"*, lead-
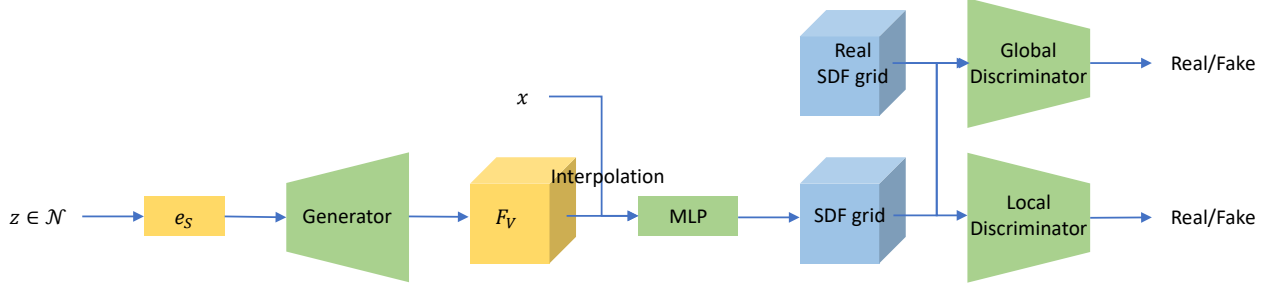
Figure 7. The network architecture of the 3D generator based on SDF-StyleGAN [79].



(a) Solid white background.



(b) Solid green background.
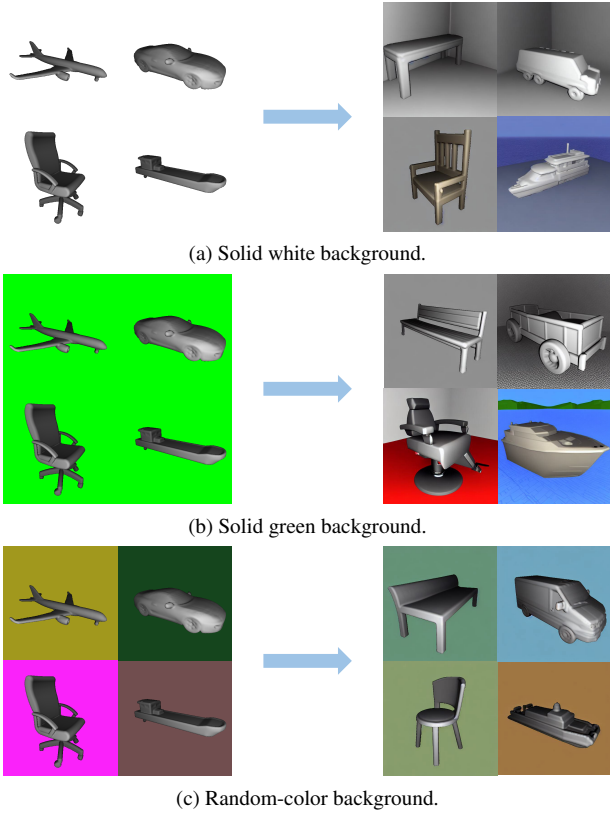


(c) Random-color background.

Figure 8. The results of fine-tuning Stable Diffusion using shape renderings with different backgrounds. For each type of background, we visualize the shape renderings used to fine-tune Stable Diffusion on the left and the images synthesized by the fine-tuned Stable Diffusion on the right.

ing to a set of 1300 generated images (denoting as $D_{gen}$). Finally, we compute the FID between $D_S$ and $D_{gen}$. As Tab. 4 shows, augmenting the backgrounds of shape renderings with random colors significantly boosts the FID, which demonstrates its effectiveness.

## D. Details of 3D Optimization with 3D Shape Prior

### D.1. DVGO-based Volume Rendering

In the optimization stage, we adopt DVGO [65] as our 3D scene representation which represents NeRF [38] with a density voxel grid $V_{density} \in \mathbb{R}^{N_x \times N_y \times N_z}$ and a shallow color MLP network $f_{rgb}$ for efficient optimization. Given a 3D position $x$, we query its density $\sigma$ and color $c$ by:

$$\tilde{\sigma} = \text{interp}(x, V_{density}), \tag{6a}$$

$$\sigma = \text{softplus}(\tilde{\sigma}) = \log(1 + \exp(\tilde{\sigma} + b)), \tag{6b}$$

$$c = f_{rgb}(\gamma(x)), \tag{6c}$$

where $\gamma(\cdot)$ denotes a positional encoding function. $\text{interp}(\cdot)$ denotes the trilinear interpolation. The shifted softplus function $\text{softplus}(\cdot)$ is applied to transform the raw density value $\tilde{\sigma}$ into activated density value $\sigma$ (i.e., a mapping of $\mathbb{R} \to \mathbb{R}_{\geq 0}$), the shift $b$ is a hyperparameter. To be noted, the density grid $V_{density}$ stores the raw density values instead of the activated ones. DVGO [65] calls the scheme of interpolating on the raw density values first and then performing softplus activation as "post-activation" and demonstrates its advantages on producing sharper shape boundaries over other choices.

To render the color of a pixel $\hat{C}(r)$, we cast the ray $r$ from the camera center through the pixel, and sample $K$ points between the pre-defined near and far planes. We then query the densities and colors of the $K$ ordered sampled points $\{(\sigma_i, c_i)\}_{i=1}^K$ using Eq. (6). Finally, we accumulate the $K$ queried results into a single color with the volume rendering process:

$$\hat{C}(r) = \left(\sum_{i=1}^K T_i \alpha_i c_i\right) + T_{K+1} c_{bg}, \tag{7a}$$

$$\alpha_i = \text{alpha}(\sigma_i, \delta_i) = 1 - \exp(-\sigma_i \delta_i), \tag{7b}$$

$$T_i = \prod_{j=1}^{i-1}(1 - \alpha_j), \tag{7c}$$
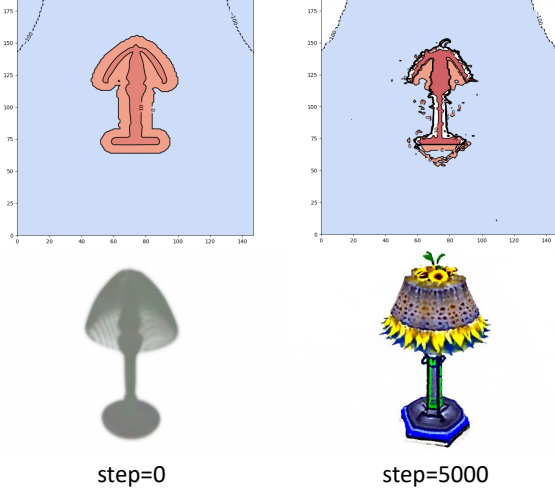
14

step=0          step=5000

Figure 9. Illustration on 3D optimization with 3D shape prior. The text prompt is *"A lamp imitating sunflower"*. We visualize the contours of the density grid (up) and the volume-rendered images (bottom) at the $0^{th}$ and $5000^{th}$ optimization steps to show how the density grid is initialized and optimized.

where $\alpha_i$ denotes the opacity representing the probability of termination at point $i$, $T_i$ denotes the accumulated transmittance from the near plane to point $i$, $\delta_i$ denotes the distance to the adjacent sampled point, and $c_{bg}$ demotes a predefined background color.

Following DVGO, all values in $\boldsymbol{V}_{density}$ are initialized as 0 and the bias term in Eq. (6b) is set to

$$b = \log\left((1 - \alpha_{init})^{-\frac{1}{s}} - 1\right), \qquad (8)$$

where $\alpha_{init}$ is a hyperparameter and is set to $10^{-6}$ in practice. With such an initialization, the accumulated transmittance $T_i$ is decayed by $1 - \alpha_{init} \approx 1$ for a ray that traces forward a distance of a voxel size $s$, making the scene "transparent" at the beginning of optimization.

### D.2. Shape Prior Initialization and Optimization

A big difference between our text-guided 3D synthesis framework and previous methods [25, 28, 31] is that we use an explicit *3D shape prior* to initialize the CLIP-guided optimization process, instead of optimizing from a *randomly-initialized* 3D representation. Given a 3D shape prior $S$ represented by an SDF grid $\tilde{\boldsymbol{V}}_{sdf} \in \mathbb{R}^{N_x \times N_y \times N_z}$, we use it to initialize the density voxel grid $\boldsymbol{V}_{density}$ with the following equations [44, 65, 76]:

$$\boldsymbol{\Sigma} = \frac{1}{\beta} \operatorname{sigmoid}\left(-\frac{\tilde{\boldsymbol{V}}_{sdf}}{\beta}\right), \qquad (9a)$$

$$\boldsymbol{V}_{density} = \max(0, \operatorname{softplus}^{-1}(\boldsymbol{\Sigma})), \qquad (9b)$$

| Method | FID $\downarrow$ | FPD $\downarrow$ | MMD $\uparrow$ |
|---|---|---|---|
| CLIP-Forge [59] | 112.38 | 6.896 | 0.670 |
| Ours | **40.83** | **1.301** | **0.725** |

Table 5. Additional quantitative results compared with CLIP-Forge on text-guided shape generation.

where $\operatorname{sigmoid}(x) = 1/(1 + e^{-x})$ and $\operatorname{softplus}^{-1}(x) = \log(e^x - 1)$. Eq. (9a) converts SDF values to activated density values (equivalent to the $\sigma$ in Eq. (6b)), where $\beta > 0$ controls the sharpness of the shape boundary, and smaller $\beta$ leads to a sharper shape boundary. We set $\beta = 0.05$ in our experiments. Eq. (9b) further transforms the activated density values into raw density values (equivalent to the $\tilde{\sigma}$ in Eq. (6a)).

With such an initialization, the density values on the shape surface will be close to $\frac{1}{2\beta} (\log(\exp(\frac{1}{\beta} \cdot \operatorname{sigmoid}(0)) - 1) \approx \frac{1}{2\beta})$. The area inside the shape surface will have larger density values ($> \frac{1}{2\beta}$), and the density values outside the shape will decrease with the distance from the shape surface. We set the minimum density value outside the shape to 0 so the area far from the shape surface has the same initialization as the original DVGO.

As Fig. 9 shows, at the beginning of the 3D optimization process (step=0), the 3D shape prior is visible due to the larger density values around the shape surface. As a result, the area around the shape surface will dominate the volume rendering, and the density/color values in this area will be updated faster than the area far from the surface. Based on the initialization, Then subsequent CLIP-guided optimization process further provides more flexibility and is able to synthesize more diverse structures and textures.

### E. Additional Results on Text-to-Shape Generation

We show additional qualitative text-guided 3D shape generation results in Fig. 10. Compared to CLIP-Forge [59], our method produces more plausible 3D shapes thanks to the high-quality 3D generator, while the shapes generated by [59] suffer from rough surfaces and discontinuities.

Besides, we also provide more quantitative comparisons with CLIP-Forge on text-to-shape generation. We generate 3 shapes for each text prompt in the text prompt set provided by CLIP-Forge and measure three metrics: 1) Fréchet Inception Distance (FID) [22] between 5 rendered images for each shape with different camera poses and a set of images rendered from the ground truth shapes in the ShapeNet dataset with the same camera poses. 2) Fréchet Point Distance (FPD) [63], for each generated shape and each ground truth shape in the ShapeNet test set, we extract the mesh at $64^3$ resolution and sample 2048 points from the mesh sur-
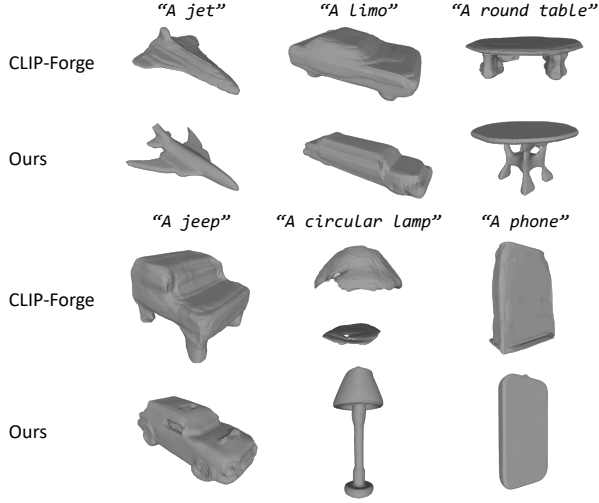
Figure 10. Additional text-guided shape generation results. All meshes are extracted at $64^3$ resolution.

face, then pass the points to a DGCNN [68] backbone network pre-trained on the point cloud classification task and use the feature of the last layer to compute this metric. 3) Maximum Measure Distance (MMD), for each generated shape represented by a $32^3$ occupancy grid, we match a shape in the ShapeNet test set based on the highest IOU, and then average the IOU across all the text queries. As Tab. 5 shows, our text-to-shape generation method outperforms CLIP-Forge on all three metrics.

## F. Additional Results on Text-to-3D Synthesis

In this section, we show additional qualitative comparison results on text-to-3D synthesis with baseline methods in Fig. 11 and more diversified generation results of our method in Fig. 12. It can be seen that our method can synthesize plausible 3D structures with the help of 3D shape priors. To better visualize the 3D structures generated by different methods, we also show video examples in the attached MP4 file.

Figure 11. Additional qualitative comparisons on text-guided 3D synthesis. For each of our results (the last row), we also visualize the 3D shape prior used to initialize the CLIP-guided optimization process below it.
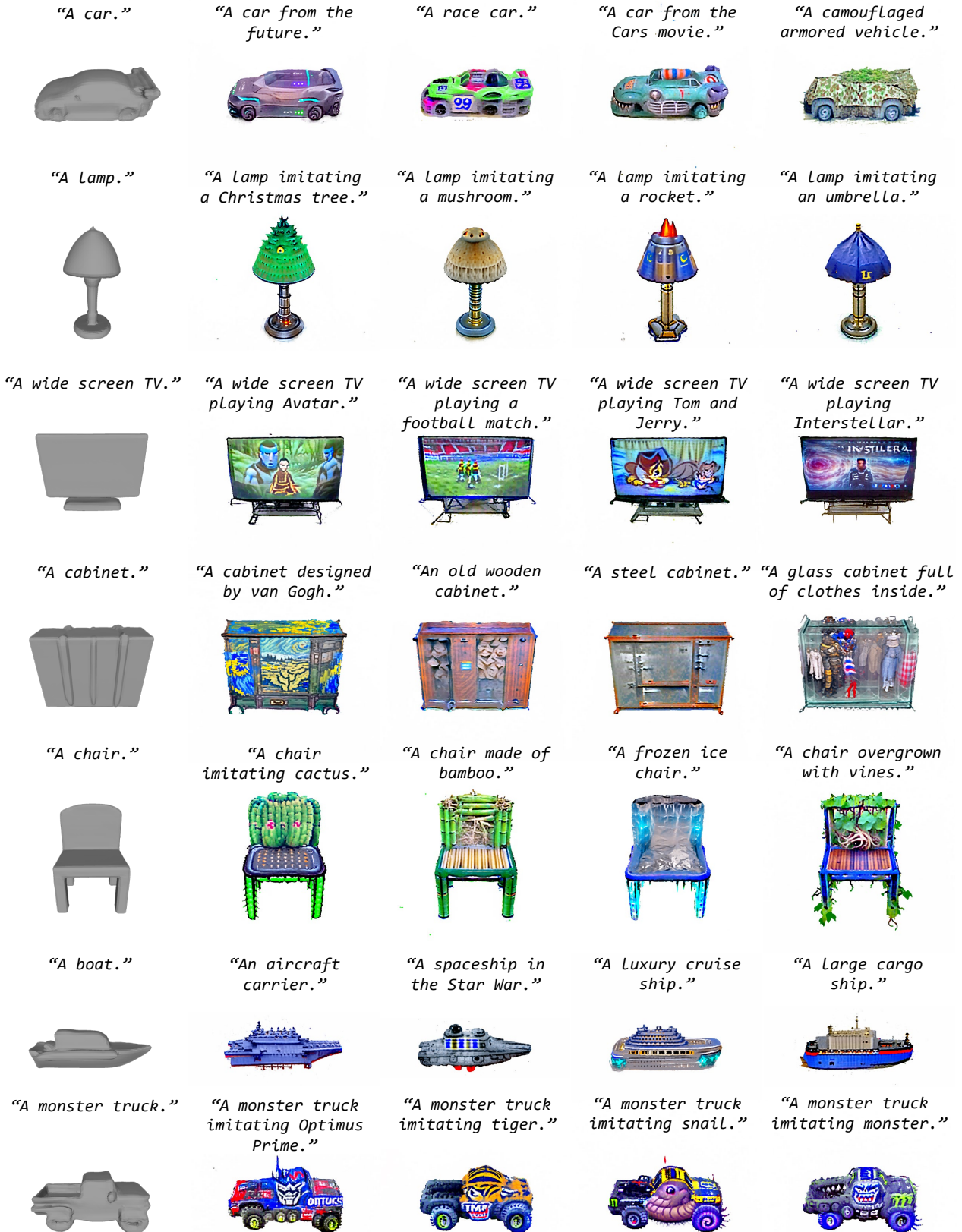
"A car."  "A car from the future."  "A race car."  "A car from the Cars movie."  "A camouflaged armored vehicle."

"A lamp."  "A lamp imitating a Christmas tree."  "A lamp imitating a mushroom."  "A lamp imitating a rocket."  "A lamp imitating an umbrella."

"A wide screen TV."  "A wide screen TV playing Avatar."  "A wide screen TV playing a football match."  "A wide screen TV playing Tom and Jerry."  "A wide screen TV playing Interstellar."

"A cabinet."  "A cabinet designed by van Gogh."  "An old wooden cabinet."  "A steel cabinet."  "A glass cabinet full of clothes inside."

"A chair."  "A chair imitating cactus."  "A chair made of bamboo."  "A frozen ice chair."  "A chair overgrown with vines."

"A boat."  "An aircraft carrier."  "A spaceship in the Star War."  "A luxury cruise ship."  "A large cargo ship."

"A monster truck."  "A monster truck imitating Optimus Prime."  "A monster truck imitating tiger."  "A monster truck imitating snail."  "A monster truck imitating monster."

Figure 12. Additional text-to-3D synthesis results. We visualize the 3D shape prior used for optimization in the first column.

## G. Integration with SVR Models

### G.1. Text-to-Shape Generation using SVR models

Single-view reconstruction (SVR) models can reconstruct a 3D shape from a single input image. We then ask, can we use an SVR model directly as the image-to-shape module in our framework? To answer this question, we conduct the same fine-tuning process to fine-tune a Stable Diffusion model with the ShapeNet renderings provided by Choy *et al*. [9] which are commonly used by many SVR methods. We find that although the shape renderings in Choy *et al*. [9] have more complex textures, the fine-tuned model can still capture the style successfully and synthesize novel images imitating the style. With such a fine-tuned Stable Diffusion, we can solve the text-to-shape generation in a precise way: synthesize an image using the fine-tuned Stable Diffusion with text prompt in the format of *"a CLS in the style of *"*, and then directly feed the synthesized image into the SVR model. We show some text-guided shape generation results using two SVR methods, *i.e*., occupancy networks [37] and DVR [42], in Fig. 13 and Fig. 14, respectively. Both methods are trained with the shape renderings provided by Choy *et al*. [9]. The occupancy networks only predict shape, while DVR can predict both shape and color. As Fig. 13 and Fig. 14 show, we achieve text-to-shape generation successfully with the synthesized images, which proves the strong generation ability of Stable Diffusion and the effectiveness of our fine-tuning pipeline.

A recent work named ISS [33] also utilizes an SVR model to perform text-to-shape generation. However, the pipeline of ISS is much more complicated. It trains a mapper network to map CLIP features to the latent space of the SVR model, which requires a two-stage fine-tuning to align the text and shape feature spaces. At inference time, ISS needs to fine-tune the mapper network for each text prompt, which is redundant in our pipeline. With the help of the fine-tuned Stable Diffusion, we can directly generate an image from the text prompt and feed the image into the SVR model to synthesize a 3D shape. Besides, thanks to the strong generation ability of Stable Diffusion, we can enjoy a much larger generation diversity and synthesize as many 3D shapes as we want for each text prompt.

### G.2. Text-to-3D Synthesis using SVR models

Despite the success in text-guided shape generation with SVR models, we find that current SVR models are very sensitive to the input images. Although we can successfully capture the style of the shape renderings using the fine-tuned Stable Diffusion, some minor flaws in the synthesized images such as offsets of the objects from the image center and unrealistic artifacts (*e.g*., a chair lacks a leg) are inevitable. These minor flaws may lead to failed shape reconstructions, whose quality affects 3D shape priors. This



"A boat with sail"  "A children chair with little legs"

"A swivel chair"  "A wooden table"

"A mushroom-like lamp"  "A sofa with legs"

Figure 13. Text-guided shape generation using fine-tuned Stable Diffusion and Occupancy Networks [37]. For each text prompt, we visualize the image synthesized by the fine-tuned Stable Diffusion on the left and the reconstructed shape on the right.



"A green SUV"  "A computer monitor"

"A red recliner seems comfortable"  "A round shaped single legged wooden table"

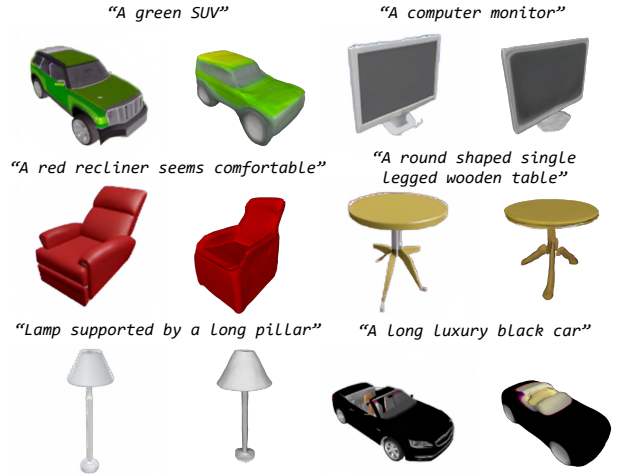"Lamp supported by a long pillar"  "A long luxury black car"

Figure 14. Text-guided shape generation using fine-tuned Stable Diffusion and DVR [42]. For each text prompt, we visualize the image synthesized by the fine-tuned Stable Diffusion on the left and the reconstructed shape on the right.

sensitiveness makes the 3D prior generation in the first stage of our framework unstable. Therefore, we choose to use a 3D generator associated with a shape embedding mapping network to generate 3D shapes in the latent shape embedding space, instead of directly using an SVR model in our framework.

We visualize six text-to-3D synthesis results using 3D shape priors produced by the occupancy networks [37] in Fig. 15. The successful results in the first two rows show the probability of integrating as SVR model into our framework. In the last row, we show two failure cases in which the SVR model fails to reconstruct plausible 3D shape pri-
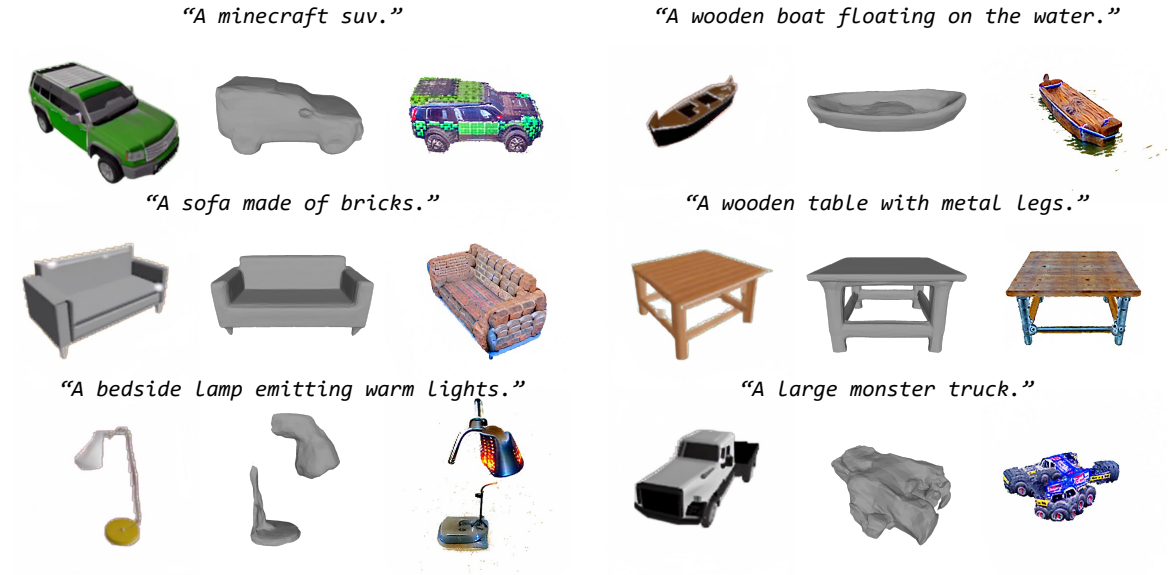
Figure 15. Text-to-3D synthesis results using occupancy networks. For each text prompt, we visualize the shape rendering image synthesized by the fine-tuned Stable Diffusion on the left, the shape reconstructed by the SVR model in the middle, and the optimization result on the right. The last row shows two failure cases.

ors to illustrate the drawbacks of using SVR models. We can observe that the discontinuity in the "bedside lamp" shape leads to discontinuity in the final optimization result, while the failed truck shape results in total chaos.