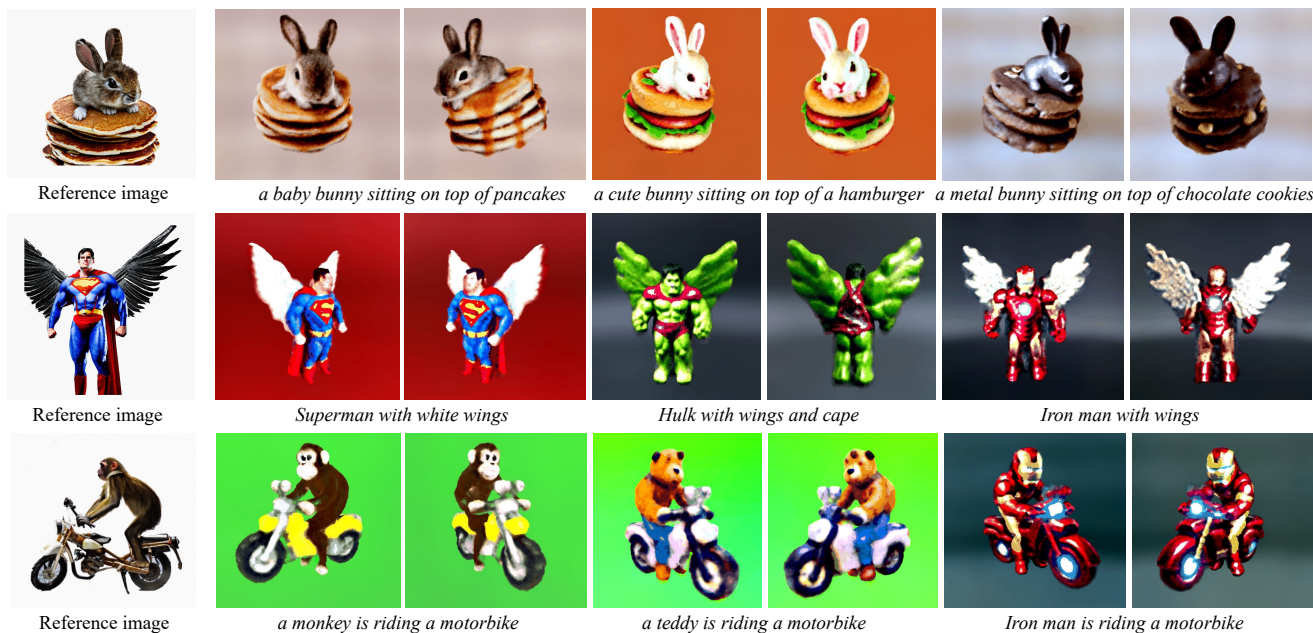


# Points-to-3D: Bridging the Gap between Sparse Points and Shape-Controllable Text-to-3D Generation

Chaohui Yu, Qiang Zhou, Jingliang Li, Zhe Zhang, Zhibin Wang, Fan Wang  
DAMO Academy, Alibaba Group



**Figure 1: Points-to-3D can create flexible 3D content with a similar shape to a single reference image. The provided reference image can be a real image or a synthesized image generated by text-to-image diffusion models, e.g., Stable Diffusion.**

## ABSTRACT

Text-to-3D generation has recently garnered significant attention, fueled by 2D diffusion models trained on billions of image-text pairs. Existing methods primarily rely on score distillation to leverage the 2D diffusion priors to supervise the generation of 3D models, e.g., NeRF. However, score distillation is prone to suffer the view inconsistency problem, and implicit NeRF modeling can also lead to an arbitrary shape, thus leading to less realistic and uncontrollable 3D generation. In this work, we propose a flexible framework of Points-to-3D to bridge the gap between sparse yet freely available 3D points and realistic shape-controllable 3D generation by distilling the knowledge from both 2D and 3D diffusion models. The core idea of Points-to-3D is to introduce controllable sparse 3D points to guide the text-to-3D generation. Specifically, we use the sparse point cloud generated from the 3D diffusion model, Point-E, as the geometric prior, conditioned on a single reference image. To better

utilize the sparse 3D points, we propose an efficient point cloud guidance loss to adaptively drive the NeRF’s geometry to align with the shape of the sparse 3D points. In addition to controlling the geometry, we propose to optimize the NeRF for a more view-consistent appearance. To be specific, we perform score distillation to the publicly available 2D image diffusion model ControlNet, conditioned on text as well as depth map of the learned compact geometry. Qualitative and quantitative comparisons demonstrate that Points-to-3D improves view consistency and achieves good shape controllability for text-to-3D generation. Points-to-3D provides users with a new way to improve and control text-to-3D generation.

## CCS CONCEPTS

• **Computing methodologies** → *Visibility; Appearance and texture representations.*

## KEYWORDS

text-to-3D, diffusion models, NeRF, point cloud

## ACM Reference Format:

Chaohui Yu, Qiang Zhou, Jingliang Li, Zhe Zhang, Zhibin Wang, Fan Wang, DAMO Academy, Alibaba Group, . 2018. Points-to-3D: Bridging the Gap between Sparse Points and Shape-Controllable Text-to-3D Generation. In *Proceedings of Make sure to enter the correct conference title from your rights*

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

Conference ACM MM '23, Oct. 29–Nov. 03, 2023, Ottawa, Canada

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

confirmation email (Conference ACM MM '23). ACM, New York, NY, USA, 10 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Recently, phenomenal advancements have been made in the field of text-to-image generation [38, 39, 42, 44, 59], mainly due to the significant achievements in large aligned image-text datasets [47], vision-language pre-training models [20, 24, 37], and diffusion models [9, 15, 42]. Inspired by these text-to-image generation results, many works have explored text-conditional diffusion models in other modalities, *e.g.*, text-to-video [16, 17, 48] and text-to-3D [19, 25, 28, 36, 54]. In this work, we focus specifically on the field of text-to-3D generation, which aims to create 3D content and can potentially be applied to many applications, *e.g.*, gaming, virtual or augmented reality, and robotic applications.

Training text-to-3D generative models can be challenging since it is difficult to attain plentiful text and 3D data pairs compared to 2D images. Most recently, DreamFusion [36] first addresses the challenge by using score distillation from a pre-trained 2D text-to-image diffusion model [44] to optimize a Neural Radiance Fields (NeRF) [29] to perform text-to-3D synthesis. The following literatures [28, 54] also use the score distillation paradigm. These methods provide and verify the solution for text-to-3D content generation without requiring 3D supervision. Despite their considerable promise, these methods are plagued by a notable issue known as the multi-face problem, or *Janus problem*, which results in inconsistencies across views. Besides, another important issue in text-to-3D generation is the lack of control over the shape of the generated 3D objects, *i.e.*, these methods may produce objects with arbitrary shapes that meet the requirements of the input text by setting different seeds. Latent-NeRF [28] first introduces sketch-shape guided 3D generation, which uses a predefined mesh as a target to supervise the geometry learning of the NeRF. However, this approach is costly and time-consuming, as it requires the predefinition of a mesh shape for each 3D generation every time.

This has motivated us to explore the possibility of cultivating prior knowledge in both 2D and 3D diffusion models to guide both the appearance and geometry learning of text-to-3D generation. Inspired by the conditional control paradigm in text-to-image diffusion models, *e.g.*, ControlNet [59] and T2I-Adapter [32], which use extra conditions (*e.g.*, sketch, mask, depth) with text prompts to guide the generation process, achieving more controllability and spatial consistency of the image. We seek a way to incorporate this conditional control mechanism into text-to-3D generation.

In this work, we propose a novel and flexible framework, dubbed Points-to-3D, to improve view consistency across views and achieve flexible controllability over 3D shapes for text-to-3D generation. The core idea of Points-to-3D is to introduce controllable sparse 3D points to guide the text-to-3D generation in terms of geometry and appearance. To achieve this, inspired by Point-E [35], we propose to distill the sparse point clouds from pre-trained 3D point cloud diffusion models as the geometry prior. These sparse 3D points are conditioned on a single reference image, which can be provided either by the user or generated by a text-to-image model. However, it is not trivial to leverage the generated sparse point clouds, which only contain 4096 3D points. To overcome this issue, we propose

a point cloud guidance loss to encourage the geometry of a randomly initialized NeRF to closely resemble the shape depicted in the reference image.

In addition to geometry, we propose to optimize the appearance conditioned on text prompt as well as the learned depth map. More concretely, we perform score distillation [28, 36] to the publicly available and more controllable 2D image diffusion models, ControlNet [59], in a compact latent space. Our approach, Points-to-3D, can bridge the gap between sparse 3D points and realistic shape-controllable 3D generation by distilling the knowledge of 2D and 3D diffusion priors. As depicted in Figure 1, given an imaginative reference image, Points-to-3D can generate realistic and shape-controllable 3D contents that vary with different text prompts.

In summary, the contributions of this paper are as follows:

- We present a novel and flexible text-to-3D generation framework, named Points-to-3D, which bridges the gap between sparse 3D points and more realistic and shape-controllable 3D generation by distilling the knowledge from pre-trained 2D and 3D diffusion models.
- To take full advantage of the sparse 3D points, we propose an efficient point cloud guidance loss to optimize the geometry of NeRF, and learn geometry-consistent appearance via score distillation by using ControlNet conditioned on text and learned depth map.
- Experimental results show that Points-to-3D can significantly alleviate inconsistency across views and achieve good controllability over 3D shapes for text-to-3D generation.

## 2 RELATED WORK

**Text-to-Image Generation.** Image generation achieves the first breakthrough results when encountering Generative Adversarial Networks (GANs) [13, 21], which train a generator to synthesize images that are indistinguishable from real images. Recently, image generation has achieved another phenomenal progress with the development of diffusion models [49]. With the improvements in modeling [9, 15], denoising diffusion models can generate various high-quality images by iteratively denoising a noised image. In addition to image-driven unconditional generative, diffusion models can generate text-conditioned images from text descriptions [38, 44]. The following works propose to add more conditions to text-to-image generation, including semantic segmentation [42], reference images [43], sketch [53], depth map [32, 59], and other conditions [18, 32, 59], which greatly promote the development and application of text-to-image generation. Driven by the success of text-to-image diffusion models, many works have explored text-conditional diffusion models in other modalities, *e.g.*, text-based manipulation [4], text-to-video [17, 48], and text-to-3D [25, 28, 36, 54]. In this work, we focus on the field of text-to-3D generation.

**Neural Radiance Fields (NeRF).** There is plenty of work on 3D scene representation, including 3D voxel grids [51], mesh [11], point clouds [1, 27, 30, 60], and implicit NeRF [29, 34]. In recent years, as a series of inverse rendering methods, NeRF-based methods have emerged as an important technique in 3D scene representation, which are capable of synthesizing novel views and reconstructing geometry surface [29, 34, 55]. Specifically, NeRF [29] represents scenes as density and radiance fields with the neural network

(MLP), allowing for photorealistic novel view synthesis. However, the computational cost of densely querying the neural network in 3D space is substantial. To improve the efficiency of NeRF, recent research has explored designing hybrid or explicit structures based on NeRF [6, 34, 51] to achieve fast convergence for radiance field reconstruction, as well as accelerating the rendering speed of NeRF [12, 14, 40]. Most of these methods require multiple views and corresponding camera parameters for training, which can not be always satisfied, especially in novel text-to-3D content generation. In this work, we view NeRF as a basic scene representation model and focus on devising a new framework for text-to-3D generation.

**Single Image 3D Reconstruction.** Various approaches exist for single image 3D reconstruction, which aims at reconstructing the object present in the image. Different formats can be used to represent the reconstructed object, such as voxels [7, 57], polygonal meshes [56], point clouds [10], and more recently, NeRFs [33, 58]. However, these methods are typically trained and evaluated on specific 3D datasets [5], making generalization to general 3D reconstruction challenging due to the lack of sufficient 3D training data. Recently, Point-E [35] explores an efficient method for general 3D content generation in the form of point clouds. It first generates a single synthetic image using a pre-trained text-to-image diffusion model, and then produces a sparse (4096 points) 3D point cloud using a point cloud diffusion model, which is conditioned on the generated image. The generalization ability of Point-E is attributed to its training on several millions of 3D data [35]. In this work, we innovatively leverage Point-E as a point cloud foundation model, to provide sparse geometry guidance for more realistic and shape-controllable text-to-3D generation.

**Text-to-3D Generation.** In recent times, the progress in text-to-image generation and 3D scene modeling has sparked a growing interest in text-to-3D content generation. Earlier work like CLIP-forge [45] consists of an implicit autoencoder conditioned on shape codes and a normalizing flow model to sample shape embeddings from textual input. However, it needs 3D training data in voxel representation, which is difficult to scale in real applications. Pure-CLIPNeRF [23] uses pre-trained CLIP [37] for guidance with a voxel grid model for scene representation to perform text-to-3D generation without access to any 3D datasets. CLIP-Mesh [31] presents a method for zero-shot 3D generation using a textual prompt, it also relies on a pre-trained CLIP model that compares the input text with differentially rendered images of the generated 3D model. DreamFields [19] first proposes to optimize the 3D representation of NeRF [29], by employing a pre-trained CLIP as guidance as well, such that all rendering views of NeRF are encouraged to match the text prompt.

More recently, DreamFusion [36] proposes to utilize a powerful pre-trained 2D text-to-image diffusion model [44] to perform text-to-3D synthesis. They propose a Score Distillation Sampling (SDS) loss to supervise the rendered views of 3D objects modeled by NeRF. The following Stable-DreamFusion [52], Latent-NeRF [28], and SJC [54] adapt the score distillation to the publicly available and computationally efficient Stable Diffusion model [42], which apply the diffusion process in a compact latent space and facilitate the development of text-to-3D generation. We build upon these works and propose a flexible Points-to-3D framework for text-to-3D

generation by bridging the gap between sparse 3D points and more realistic shape-controllable 3D content generation.

## 3 APPROACH

### 3.1 Preliminaries

In this section, we provide a brief introduction to some of the key concepts that are necessary for understanding our proposed framework in Section 3.2.

**Diffusion Model.** Diffusion models are first proposed by [49] and recently promoted by [15, 50]. A diffusion model usually consists of a forward process  $q$  that gradually adds noise to the image  $x \in X$ , and a reverse process  $p$  of gradually removing noise from the noisy data. The forward process  $q$  can be formulated as:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}), \quad (1)$$

where timesteps  $t \in [0, T]$ ,  $\beta_t$  denotes noise schedule. DDPM [15] proposes to directly attain a given timestep of the noising procedure:

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \quad (2)$$

where  $\bar{\alpha}_t = \prod_0^t 1 - \beta_t$ , and  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ . The denoising process  $p_\theta(x_{t-1}|x_t)$  starts from random noise and slowly reverses the noising process. DDPM [15] proposes to parameterize the distribution by modeling the added noise  $\epsilon$ . Recently, latent diffusion model (LDM), as a specific form of diffusion model, has achieved great progress in text-to-image generation. The well-known Stable Diffusion [42] and ControlNet [59] are both latent diffusion models.

**Score Distillation Sampling (SDS).** Score distillation sampling (SDS) is first proposed by DreamFusion [36], which achieves text-to-3D creation by incorporating two modules: a scene representation model [3] and a pre-trained text-to-image diffusion model (Imagen [44]). During training, a learnable NeRF model  $\theta$  first performs view synthesizes with a differentiable render  $g: x = g(\theta)$ , which can render an image at a given random camera pose. Then, random noise is added to  $x$  and the diffusion model  $\phi$  is to predict the added noise  $\epsilon$  from the noisy image with a learned denoising function  $\epsilon_\phi(x_t; y, t)$  given the noisy image  $x_t$ , text embedding  $y$ , and noise level  $t$ . This score function provides gradient to update the NeRF parameters  $\theta$ , which is calculated as:

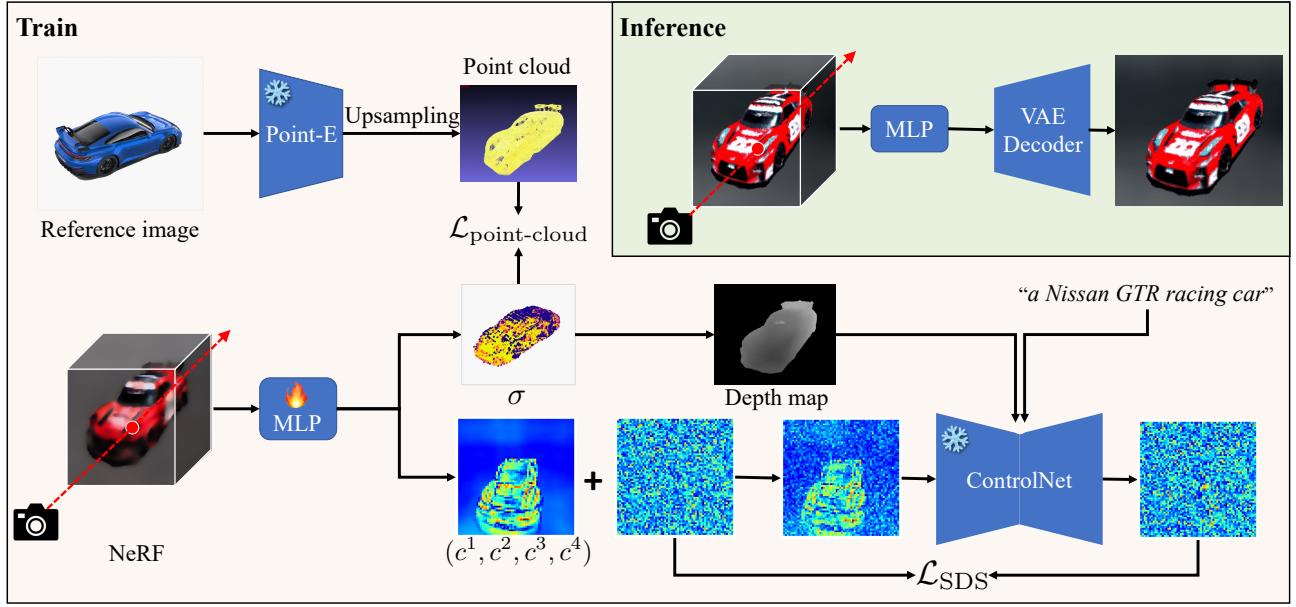
$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon} \left[ \omega(t) (\epsilon_\phi(x_t; y, t) - \epsilon) \frac{\partial x}{\partial \theta} \right], \quad (3)$$

where  $\omega(t)$  is a weighting function that depends on  $\alpha_t$ . Inspired by Stable-DreamFusion [52] and Latent-NeRF [28], which use Stable Diffusion [42], we propose to perform score distillation with a more controllable LDM, ControlNet [59], to generate more realistic and shape-controllable 3D contents.

### 3.2 Points-to-3D

In this section, we elaborate on our Points-to-3D framework, which is depicted in Figure 2.

**Architecture.** First of all, we describe the architecture of our Points-to-3D framework. As shown in Figure 2, Points-to-3D mainly consists of three models: a scene representation model (a coordinate-based MLP [34]), a text-to-image 2D diffusion model (ControlNet [59]), and a point cloud 3D diffusion model (Point-E [35]).



**Figure 2: Illustration of the proposed Points-to-3D framework for text-to-3D generation. Points-to-3D mainly consists of three parts: a scene representation model (a coordinate-based NeRF [34]), a text-to-image 2D diffusion model (ControlNet [59]), and a point cloud 3D diffusion model (Point-E [35]). During training, both 2D and 3D diffusion models are frozen.**

- **Scene Model.** Neural Radiance Field (NeRF) [29] has been an important technique used for scene representation, comprising of a volumetric raytracer and an MLP. Previous literature [28, 36, 54] has used NeRF as the scene representation model for text-to-3D generation, mainly because a NeRF model can implicitly impose the spatial consistency between different views owing to the spatial radiance field and rendering paradigm. A NeRF model usually produces a volumetric density  $\sigma$  and an RGB color  $c$ . In this work, we adopt the efficient design of Latent-NeRF [28] that produces five outputs, including the volume density  $\sigma$  and four pseudo-color channels  $\{C = (c^1, c^2, c^3, c^4)\} \in \mathbb{R}^{64 \times 64 \times 4}$  that correspond to the four input latent features for latent diffusion models [42]:

$$(c^1, c^2, c^3, c^4, \sigma) = \text{MLP}(x, y, z, d; \theta), \quad (4)$$

where  $x, y, z$  denote 3D coordinates,  $d$  is the view direction. We use Instant-NGP [34] as the scene representation model by default.

- **Text-to-Image 2D Diffusion Model.** Since Imagen [44] used by DreamFusion [36] is not publicly available, we use Stable Diffusion as the text-to-image diffusion model initially, as previously explored in existing literature [28, 52, 54]. However, the original Stable Diffusion v1.5 is not controllable to support additional input conditions. In this work, we first propose to use the pre-trained ControlNet [59] conditioned on depth map as the 2D diffusion model in Points-to-3D. As depicted in Figure 2, in addition to the input text prompt, e.g., “a Nissan GTR racing car”, we further introduce the predicted depth map  $M \in \mathbb{R}^{H \times W \times 1}$  of our NeRF model as the conditional control. The depth map is computed as follows, for simplicity, we only show the depth value calculation on one pixel:

$$M_i = \sum_{k=1}^K w_k t_k, \quad (5)$$

and

$$w_k = \alpha_k \prod_{j < k} (1 - \alpha_j), \text{ and } \alpha_k = 1 - \exp(-\sigma_k \|t_k - t_{k+1}\|). \quad (6)$$

where  $K$  is the total number of sampling points along a ray, and  $t_k$  denotes the depth hypothesis at point  $k$ . The better and more accurate the predicted depth map  $M$ , the more geometrically consistent views ControlNet will synthesize.

- **Point Cloud 3D Diffusion Model.** To control the geometry of NeRF for text-to-3D generation, we propose in this paper, for the first time, the distillation of prior knowledge from the pre-trained large point cloud diffusion model, Point-E [35]. Point-E [35] is an efficient 3D diffusion model for generating sparse 3D point clouds (4096 points) from text prompts or images in about 1 minute. As illustrated in Figure 2, we utilize the pre-trained Point-E model to regulate the geometry learning of NeRF. Specifically, the model generates a sparse 3D point cloud consisting of 4096 points, which is conditioned on a reference image and can flexibly represent the object’s shape depicted in the image. However, it is not trivial to guide the NeRF’s geometry with only sparse 3D points, we propose a sparse point cloud guidance loss  $\mathcal{L}_{\text{point-cloud}}$  to address this issue, which is illustrated in the next section.

It is worth noting that Points-to-3D enables users to easily control the shape of the generated content by providing a reference image, which can be a real image or a generated image via text-to-image models [32, 42, 59].

**Sparse 3D Points Guidance.** The core idea of our Points-to-3D is to introduce controllable sparse 3D points to guide the text-to-3D generation. In this section, we elaborate on how to leverage the sparse 3D points. It is challenging to use a sparse 3D point cloud to guide the geometry learning of NeRF. Previous work on improving



NeRF’s geometry uses the depth of sparse points to supervise the predicted depth [8, 41]. However, the 3D points are computed using multiple views via COLMAP [46], and the information about which view each 3D point belongs to has been calculated in advance. In our case, only a single RGB image is used to generate the sparse 3D points, when we project all the points to the current view to attain a sparse depth map, there will be aliasing problems between the front and the rear 3D points.

In this work, we present a sparse point cloud guidance loss. Specifically, let  $P_s = \{(x_i, y_i, z_i)\}_{i=1}^{4096}$  be the original sparse 3D points generated by Point-E [35] conditioned on a reference image. Instead of using  $P_s$  directly, we experimentally find that making the sparse point cloud to be dense can provide better geometry supervision and produce more realistic 3D contents. We propose to upsample  $P_s$  by iteratively performing 3D points interpolation via a simple rule, *i.e.*, for each point  $p_i$ , we add a new 3D point at the middle position between each of its nearest  $q$  neighbor points and  $p_i$ . The process is depicted in Figure 3. We set  $q = 20, n = 2$  by default. Now we get the dense 3D points  $P_d$ , which contain about 500k points after eliminating duplicate points.

Ideally, we want to align the geometry (the volume density  $\sigma$ ) of NeRF with the shape of  $P_d$  to ensure that the generated 3D content of Points-to-3D closely resembles the reference image. In addition, we also want to provide NeRF with a level of flexibility and adaptability in its geometry to enable the generation of new details while satisfying different text prompts. Instead of using the per-view sparse depth map supervision, which has a front-rear aliasing issue as discussed above, and is also not efficient as it only optimizes the current view’s depth, we propose an efficient point cloud guidance loss  $\mathcal{L}_{\text{point-cloud}}$  to directly optimize the whole geometry ( $\sigma$ ) in 3D space. Specifically, we encourage the occupancy ( $\alpha$ ) corresponding to the NeRF points  $P_{\text{nerf}}$  that near the point cloud  $P_d$  to be close to 1, while the occupancy of the NeRF points that far from the point cloud  $P_d$  to be close to 0. Furthermore, we make the geometry capable of generating new details adaptively by ignoring the supervision of some parts of occupancy. We first compute the closest distance between each point in  $P_{\text{nerf}}$  and all points in  $P_d$ :  $\mathcal{D} = \text{Dist}(P_{\text{nerf}}, P_d)$ ,  $\mathcal{D} \in \mathbb{R}^{S \times 1}$ , where  $S$  denotes the number of points in  $P_{\text{nerf}}$ . Then, normalize  $\mathcal{D}$  via:  $\widehat{\mathcal{D}} = \frac{\mathcal{D}}{0.5 \cdot (\max(P_{\text{nerf}}) - \min(P_{\text{nerf}}))}$ . Finally, The calculation of  $\mathcal{L}_{\text{point-cloud}}$  can be formulated as:

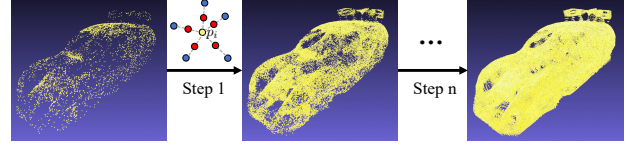
$$\mathcal{L}_{\text{point-cloud}} = \text{CrossEntropy}(\alpha(P_{\text{nerf}}), O(P_{\text{nerf}})), \quad (7)$$

and

$$O_i = \begin{cases} 1 - \widehat{\mathcal{D}}_i, & \text{if } 1 - \widehat{\mathcal{D}}_i > \tau_1; \\ 0, & \text{else if } 1 - \widehat{\mathcal{D}}_i < \tau_2; \\ -1, & \text{otherwise;} \end{cases} \quad (8)$$

where  $O(P_{\text{nerf}})$  denotes the target occupancy of all NeRF points,  $1 - \widehat{\mathcal{D}}$  indicates the degree of proximity to the guided point cloud  $P_d$ , and  $\tau_1, \tau_2$  are two hyperparameters that are experimentally set to 0.95 and 0.9 respectively. We ignore the supervision of points with  $\tau_2 < 1 - \widehat{\mathcal{D}} < \tau_1$ , allowing the model to adaptively add new details into the geometry to match the text prompts, as well as fix broken holes in the imperfect guided point cloud  $P_d$ .

**Training Objectives.** The training objectives of Points-to-3D consist of three parts: the point cloud guidance loss  $\mathcal{L}_{\text{point-cloud}}$ , the



**Figure 3: Illustration of the point cloud upsampling process. For each original 3D point (e.g.,  $p_i$ ), we add new 3D points (red points) between each of the nearest  $q$  neighbor points (blue points) and point  $p_i$  for each interpolation step.**

score distillation sampling loss  $\mathcal{L}_{\text{SDS}}$ , and a sparsity loss  $\mathcal{L}_{\text{sparse}}$ . The sparsity loss is suggested by [52], which can suppress floaters by regularizing the rendering weights:

$$\mathcal{L}_{\text{sparse}} = - \sum_k (w_k \log w_k + (1 - w_k) \log(1 - w_k)). \quad (9)$$

We introduce the depth map condition  $M$  calculated by Equation 5 and update the score distillation sampling loss in Equation 3 as follows:

$$\nabla_{\theta} \mathcal{L}_{\text{SDS}}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon} [\omega(t) (\epsilon_{\phi}(x_t; y, M, t) - \epsilon) \frac{\partial x}{\partial \theta}]. \quad (10)$$

The overall learning objective is computed as:

$$\mathcal{L} = \lambda_{\text{point}} \mathcal{L}_{\text{point-cloud}} + \lambda_{\text{SDS}} \mathcal{L}_{\text{SDS}} + \lambda_{\text{sparse}} \mathcal{L}_{\text{sparse}}. \quad (11)$$

## 4 EXPERIMENTS

### 4.1 Baselines

We consider three text-to-3D generation baselines: DreamFusion [36, 52], Latent-NeRF [28], and SJC [54]. Instead of using the close-sourced Imagen [44] diffusion model, both Latent-NeRF and SJC use the publicly available Stable Diffusion [42]. We mainly compare our Points-to-3D with Latent-NeRF and SJC in the experiments. We provide more results including comparisons with DreamFields [19], and DreamFusion [36] in our supplementary materials.

### 4.2 Implementation Details

We use Instant-NGP [34] as our scene model. Following the camera sampling method in [36], during training, a camera position is randomly sampled in spherical coordinates, and we also randomly enlarge the FOV when rendering with NeRF. In addition to the training in latent space shown in Figure 2, we experimentally find that further performing RGB refinement in RGB space, which is introduced in [28], can further improve the text-to-3D generation results. Our Points-to-3D takes less than 50 minutes per text prompt to complete a 3D generation on a single A100 GPU, and most of the time is spent on calculating  $\mathcal{L}_{\text{point-cloud}}$ . We train for 5000 iterations using AdamW optimizer with a learning rate of  $1e^{-3}$ . The hyperparameters of  $\lambda_{\text{point}}, \lambda_{\text{SDS}}, \lambda_{\text{sparse}}$  are set to  $5e^{-6}, 1.0, 5e^{-4}$ , respectively.

### 4.3 Ablation Studies

**Effect of Point Cloud Guidance Loss.** In this section, we evaluate the proposed point cloud guidance loss  $\mathcal{L}_{\text{point-cloud}}$ . Concretely, we evaluate Points-to-3D by eliminating the point cloud guidance. We also verify the per-view sparse depth map loss as discussed in Section 3.2. The results are shown in Figure 4. We first produce a

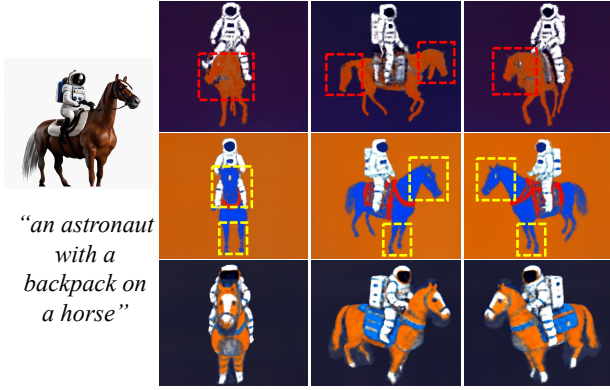


Figure 4: Illustration of the effect of our  $\mathcal{L}_{\text{point-cloud}}$ . Given a reference image and a text prompt, our Points-to-3D with  $\mathcal{L}_{\text{point-cloud}}$  (the 3rd row) can generate more realistic 3D content than both the per-view depth map loss (the 2nd row) and that without any geometry constraints [28] (the 1st row).



Figure 5: Comparison of rendered views of models trained with  $P_s$  and  $P_d$  as geometry guidance, respectively. The text prompt is “a Nissan GTR racing car”.

reference image with the text prompt: “an astronaut with a backpack on a horse” using Stable Diffusion. Then we use  $\mathcal{L}_{\text{point-cloud}}$  (the 3rd row), a designed per-view depth map loss (the 2nd row), and without any geometry constraints (the 1st row), to train three models with the same text prompt, respectively. We can find that without any geometry constraints, the generated content suffers an obvious view inconsistency problem (red dashed boxes). The result of using our designed per-view depth map loss as geometry supervision further improves the multi-face issue. However, the rendered images are less realistic and even broken (yellow dashed boxes) due to the sparsity of point clouds and the inefficiency of the per-view supervision. It is worth noting that the result of using  $\mathcal{L}_{\text{point-cloud}}$  shows more details in both “astronaut” and “horse”. That is, Points-to-3D with  $\mathcal{L}_{\text{point-cloud}}$  for geometry optimization can generate more realistic 3D content.

**Effect of 3D Points Upsampling.** In this section, we analyze the effect of upsampling the generated sparse 3D point cloud. As shown in Figure 5, we compare the rendered views of Points-to-3D trained with sparse (4096) 3D points  $P_s$  and upsampled denser (~500k) 3D points  $P_d$  as the geometry guidance, respectively. The 1st column represents the original sparse points  $P_s$  produced by Point-E [35] given the reference image shown in Figure 2, and the upsampled points  $P_d$  via our designed rule. The 2nd ~ 4th columns



Figure 6: Visualization of two 3D models trained with the same reference image (generated by Stable Diffusion [42]) and the corresponding sparse 3D points but different texts.



Figure 7: Comparison of two 3D models trained with the same reference image and sparse 3D points shown in the 1st column. The 1st and the 2nd rows denote training without and with adaptive design in  $\mathcal{L}_{\text{point-cloud}}$ , respectively. The text prompt is “a wooden chair”.

are three corresponding rendered views. We can see that the results guided by  $P_d$  are more realistic compared to those guided by  $P_s$ . This is due to that a denser point cloud can offer more supervision to encourage the NeRF to learn a more concise geometry. Moreover, better geometry (depth map) can also guide ControlNet [59] to generate more geometry-consistent and realistic images that match the input text prompt.

**Effect of Adaptive Design in  $\mathcal{L}_{\text{point-cloud}}$ .** In this section, we illustrate the effect of the adaptive design in  $\mathcal{L}_{\text{point-cloud}}$ . That is, in Equation 7 and Equation 8, we propose to ignore the supervision of those NeRF points with  $\tau_2 < 1 - \hat{D} < \tau_1$  to let Points-to-3D to adaptively adjust the geometry to match the text prompt. This adaptive design serves two main purposes: a). it offers the capacity to create new details without changing the main shape of the 3D content. b). it can fill broken holes in the imperfect point clouds  $P_d$ .

As shown in Figure 6, we visualize two generated 3D contents using Points-to-3D with the same reference image and sparse point cloud but different text prompts. The last three columns represent the rendered images, the rendered depth maps, and the rendered normals at the same camera pose, respectively. We can clearly observe that Points-to-3D can generate more specific new details to match different input text prompts based on the same point cloud guidance. In Figure 7, we analyze the effect of adaptive design in filling holes in the imperfect point cloud. Given a reference image, Point-E [35] may produce non-uniform point clouds, e.g., broken holes in the chair back in this instance. If we enforce all the NeRF points closed to the point cloud to be positive class and otherwise negative class, it is difficult to set an appropriate distance threshold for all 3D contents and will cause broken holes. For instance, we compare the results of rendered images and corresponding depth maps trained without and with adaptive design in the 1st and 2nd





**Figure 8: Qualitative comparison with Latnet-NeRF [28] and SJC [54] on single-object generation (the 1st ~ 4th rows) and scene generation (the 5th ~ 8th rows). The 1st column denotes reference images used for Points-to-3D, where the top four are real images and the bottom four are synthetic images generated using Stable Diffusion [42]. (Best viewed by zooming in.)**

row, respectively. Points-to-3D can naturally repair the broken holes in both geometry and appearance. We also analyze the effect of the depth map condition in our supplementary materials.

#### 4.4 Shape-Controllable Text-to-3D Generation

As special concepts and shapes are usually difficult to describe by text prompts but easy with images, it is desperately needed to have a mechanism to guide the text-to-3D content generation with images. In this section, we evaluate Points-to-3D in generating view-consistent and shape-controllable 3D contents with a single reference image for geometry guidance. Considering that DreamFusion [36] and Magic3D [25] use their proprietary text-to-image diffusion models [2, 44] and neither releases the code, we mainly compare with Latent-NeRF [28] and SJC [54]. As shown in Figure 8, we mainly compare two aspects: single-object generation and scene (consists of multiple objects) generation.

For the single-object generation (the 1st ~ 4th rows), Latent-NeRF [28] is easy to suffer the view inconsistency problem, and sometimes fails to generate reasonable content. SJC [54] looks a

little better than Latent-NeRF in terms of view consistency of the generated objects, however, it also sometimes fails to generate content that matches the text description (e.g., the 2nd and the 4th rows). Our Points-to-3D can automatically generate view-consistent and more realistic single objects. It is worth noting that Points-to-3D can generate more lifelike details, e.g., the logos of Converse, Nike, GUCCI, and LV.

For more challenging scene generation (the 5th ~ 8th rows), the inherent view inconsistency problem of Latent-NeRF [28] becomes more serious, e.g., multiple teapot spouts in the 6th row and multiple hands or legs in the 7th row. Besides, both Latent-NeRF and SJC can easily lose some concepts of the input text prompts, e.g., “motorbike” in the 5th row, “tray” in the 6th row, and “tuba” in the last row. In contrast, our Points-to-3D can create view-consistent 3D content and preserve the concepts contained in the text prompts.

Furthermore, Points-to-3D enables users to arbitrarily create or modify 3D content that has a similar shape to the reference image. We provide more comparisons in our supplementary materials.

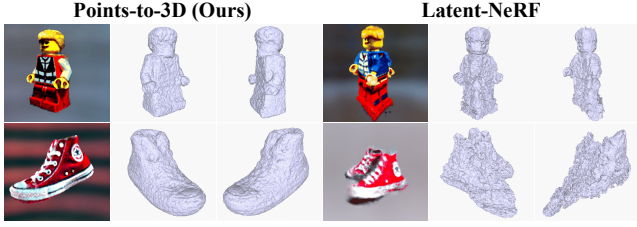


Figure 9: Mesh comparison through Marching Cubes [26].



Figure 10: Compositional generation of Points-to-3D.

#### 4.5 Geometry Comparison

We compare the learned geometry of Points-to-3D and Latent-NeRF [28], both of which use Instant-NGP [34] as the scene model. As depicted in Figure 9, we show two generation results produced using two text prompts: “a lego man” and “a red converse allstar shoe”. Each contains three views: a rendered RGB image and two views of mesh. The meshes are extracted by Marching Cubes [26] from density field of the learned Instant-NGP. We can clearly observe that compared to the flawed meshes of Latent-NeRF, Points-to-3D can generate more delicate meshes. That is, in addition to synthesis view-consistent novel views, Points-to-3D can learn controllable and more compact geometry for text-to-3D generation.

#### 4.6 Compositional Generation

We analyze the effectiveness of Points-to-3D in generating compositional 3D content. As shown in Figure 10, by taking the manually composited sparse 3D points of multiple reference images as geometry guidance, Points-to-3D can perform view-consistent and shape-controllable text-to-3D generation. The results indicate that Points-to-3D enables users to freely composite objects using multiple reference images and generate more imaginative 3D content.

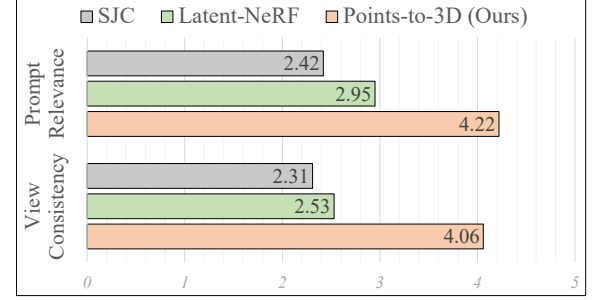
#### 4.7 Quantitative Comparisons

**CLIP R-precision.** In this section, we calculate the CLIP R-precision metrics for Latent-NeRF [28], SJC [54], and our Points-to-3D. We compute CLIP R-precision following [19] on 50 text and 3D model pairs (shown in our supplementary materials) based on three CLIP image encoders (ViT-B/16, ViT-B/32, and ViT-L/14). For each 3D generation, we randomly select two rendered views for calculation. The results are reported in Table 1, the higher scores for our Points-to-3D results indicate that renderings from our 3D model outputs more accurately resemble the text prompts.

**User Studies.** The CLIP R-precision metric focuses on the matching degree of rendered views and text prompts, but it is difficult to reflect view consistency and image realism. We conduct user studies with 22 participants to evaluate different methods based on user preferences. We ask the participants to give a preference

**Table 1: Quantitative comparison using CLIP R-precision of Latent-NeRF [28], SJC [54], and our Points-to-3D.**

Method	ViT-B/16 $\uparrow$	ViT-B/32 $\uparrow$	ViT-L/14 $\uparrow$
Latent-NeRF [28]	53.00%	59.00%	66.00%
SJC [54]	61.00%	57.00%	71.00%
Points-to-3D (Ours)	<b>81.00%</b>	<b>81.00%</b>	<b>90.00%</b>



**Figure 11: Quantitative comparison via user studies with 22 participants to measure preference in terms of view consistency and prompt relevance.**

score (range from 1 ~ 5) in terms of view consistency and prompt relevance for each anonymized method’s generation. As shown in Figure 11, we report the average scores on a randomly composed evaluation set that consists of 36 generation results of each method. We find that Points-to-3D is significantly preferred over both Latent-NeRF and SJC in terms of view consistency and prompt relevance. We provide more detailed information about the user study, please refer to our supplementary materials.

#### 5 LIMITATIONS

While Points-to-3D allows for flexible text-to-3D generation and improves over prior works in terms of realism, view consistency, and shape controllability, we observe several limitations. First, as Points-to-3D is built upon pre-trained 2D image diffusion model [59] and 3D point cloud diffusion model [35], it will be affected when ControlNet or Point-E fails with certain objects. This issue might be alleviated by developing more powerful foundation models. Second, while achieving good controllability of 3D shapes, Points-to-3D needs a single reference image for geometry guidance. This issue can be alleviated by cropping objects from real images using Segment Anything Model (SAM) [22], or direct generating an image using text-to-image models, e.g., Stable Diffusion, ControlNet.

#### 6 CONCLUSIONS

In this work, we propose Points-to-3D, a novel and flexible text-to-3D generation framework. We inspire our framework by alleviating the view inconsistency problem and improving the controllability of 3D shapes for 3D content generation. To control the learned geometry, we innovatively propose to distill the geometry knowledge (sparse 3D points) from the 3D point cloud diffusion model (Point-E). To better utilize the sparse point cloud, we propose an efficient point cloud guidance loss to adaptively align the geometry between NeRF and sparse points. Besides, to make the 3D content more realistic and view-consistent, we optimize the NeRF model conditioned on both text and the learned compact depth map,



by performing score distillation to the 2D image diffusion model (ControlNet). Both qualitative and quantitative comparisons demonstrate the superiority of Points-to-3D in generating view-consistent and shape-controllable 3D contents.

## REFERENCES

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. 2018. Learning representations and generative models for 3d point clouds. In *International conference on machine learning*. PMLR, 40–49.
- [2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. 2022. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324* (2022).
- [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. 2022. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5470–5479.
- [4] Tim Brooks, Aleksander Holynski, and Alexei A Efros. 2022. Instructpix2pix: Learning to follow image editing instructions. *arXiv preprint arXiv:2211.09800* (2022).
- [5] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. 2015. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012* (2015).
- [6] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. 2022. Tensorf: Tensorial radiance fields. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*. Springer, 333–350.
- [7] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 2016. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VIII 14*. Springer, 628–644.
- [8] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. 2022. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12882–12891.
- [9] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems* 34 (2021), 8780–8794.
- [10] Haoqiang Fan, Hao Su, and Leonidas J Guibas. 2017. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 605–613.
- [11] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. 2022. Get3d: A generative model of high quality 3d textured shapes learned from images. *Advances in Neural Information Processing Systems* 35 (2022), 31841–31854.
- [12] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. 2021. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 14346–14355.
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.
- [14] Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. 2021. Baking neural radiance fields for real-time view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5875–5884.
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.
- [16] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. 2022. Video diffusion models. *arXiv preprint arXiv:2204.03458* (2022).
- [17] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. 2022. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868* (2022).
- [18] Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. 2023. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778* (2023).
- [19] Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. 2022. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 867–876.
- [20] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*. PMLR, 4904–4916.
- [21] Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4401–4410.
- [22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. 2023. Segment Anything. *arXiv:2304.02643* (2023).
- [23] Han-Hung Lee and Angel X Chang. 2022. Understanding pure clip guidance for voxel grid nerf models. *arXiv preprint arXiv:2209.15172* (2022).
- [24] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*. PMLR, 12888–12900.
- [25] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiao-hui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. 2022. Magic3D: High-Resolution Text-to-3D Content Creation. *arXiv preprint arXiv:2211.10440* (2022).
- [26] William E Lorensen and Harvey E Cline. 1987. Marching cubes: A high resolution 3D surface construction algorithm. *ACM siggraph computer graphics* 21, 4 (1987), 163–169.
- [27] Shitong Luo and Wei Hu. 2021. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2837–2845.
- [28] Gal Metzger, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. 2022. Latent-NeRF for Shape-Guided Generation of 3D Shapes and Textures. *arXiv preprint arXiv:2211.07600* (2022).
- [29] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.
- [30] Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas J Guibas. 2019. Structrnet: Hierarchical graph networks for 3d shape generation. *arXiv preprint arXiv:1908.00575* (2019).
- [31] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Tiberiu Popa. 2022. CLIP-Mesh: Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia 2022 Conference Papers*. 1–8.
- [32] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. 2023. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453* (2023).
- [33] Norman Müller, Andrea Simonelli, Lorenzo Porzi, Samuel Rota Buló, Matthias Nießner, and Peter Kotschieder. 2022. Autof: Learning 3d object radiance fields from single view observations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3971–3980.
- [34] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)* 41, 4 (2022), 1–15.
- [35] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. 2022. Point-E: A System for Generating 3D Point Clouds from Complex Prompts. *arXiv preprint arXiv:2212.08751* (2022).
- [36] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988* (2022).
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*. PMLR, 8748–8763.
- [38] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125* (2022).
- [39] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *International Conference on Machine Learning*. PMLR, 8821–8831.
- [40] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. 2021. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 14335–14345.
- [41] Barbara Roessle, Jonathan T Barron, Ben Mildenhall, Pratul P Srinivasan, and Matthias Nießner. 2022. Dense depth priors for neural radiance fields from sparse input views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12892–12901.
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10684–10695.
- [43] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2022. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arXiv:2208.12242* (2022).
- [44] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* 35 (2022), 36479–36494.
- [45] Aditya Sanghi, Hang Chu, Joseph G Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshian. 2022. Clip-forged: Towards zero-shot text-to-shape generation. In *Proceedings of the IEEE/CVF Conference on Computer*

- Vision and Pattern Recognition*. 18603–18613.
- [46] Johannes L Schonberger and Jan-Michael Frahm. 2016. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4104–4113.
  - [47] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402* (2022).
  - [48] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. 2022. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792* (2022).
  - [49] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*. PMLR, 2256–2265.
  - [50] Yang Song and Stefano Ermon. 2019. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems* 32 (2019).
  - [51] Cheng Sun, Min Sun, and Hwann-Tzong Chen. 2022. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5459–5469.
  - [52] Jiaxiang Tang. 2022. Stable-dreamfusion: Text-to-3D with Stable-diffusion. <https://github.com/ashawkey/stable-dreamfusion>.
  - [53] Andrey Voynov, Kfir Aberman, and Daniel Cohen-Or. 2022. Sketch-Guided Text-to-Image Diffusion Models. *arXiv preprint arXiv:2211.13752* (2022).
  - [54] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. 2022. Score Jacobian Chaining: Lifting Pretrained 2D Diffusion Models for 3D Generation. *arXiv preprint arXiv:2212.00774* (2022).
  - [55] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. 2021. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689* (2021).
  - [56] Chao Wen, Yinda Zhang, Zhuwen Li, and Yanwei Fu. 2019. Pixel2mesh++: Multi-view 3d mesh generation via deformation. In *Proceedings of the IEEE/CVF international conference on computer vision*. 1042–1051.
  - [57] Haozhe Xie, Hongxun Yao, Xiaoshuai Sun, Shangchen Zhou, and Shengping Zhang. 2019. Pix2vox: Context-aware 3d reconstruction from single and multi-view images. In *Proceedings of the IEEE/CVF international conference on computer vision*. 2690–2698.
  - [58] Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Humphrey Shi, and Zhangyang Wang. 2022. Sinnerf: Training neural radiance fields on complex scenes from a single image. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII*. Springer, 736–753.
  - [59] Lvmin Zhang and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543* (2023).
  - [60] Linqi Zhou, Yilun Du, and Jiajun Wu. 2021. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5826–5835.