

# DreamLCM: Towards High Quality Text-to-3D Generation Via Latent Consistency Model

Anonymous Author(s)  
Submission Id: 1610

## ABSTRACT

Recently, the text-to-3D task has developed rapidly due to the appearance of the SDS method. However, the SDS method always generates 3D objects with poor quality due to the over-smooth issue. This issue is attributed to two factors: 1) the DDPM single-step inference produces poor guidance gradients; 2) the randomness from the input noises and timesteps averages the details of the 3D contents. In this paper, to address the issue, we propose DreamLCM which incorporates the Latent Consistency Model (LCM). DreamLCM leverages the powerful image generation capabilities inherent in LCM, enabling generating consistent and high-quality guidance, i.e., predicted noises or images. Powered by the improved guidance, the proposed method can provide accurate and detailed gradients to optimize the target 3D models. In addition, we propose two strategies to enhance the generation quality further. Firstly, we propose a guidance calibration strategy, utilizing Euler solver to calibrate the guidance distribution to accelerate 3D models to converge. Secondly, we propose a dual timestep strategy, increasing the consistency of guidance and optimizing 3D models from geometry to appearance in DreamLCM. Experiments show that DreamLCM achieves state-of-the-art results in both generation quality and training efficiency.

## CCS CONCEPTS

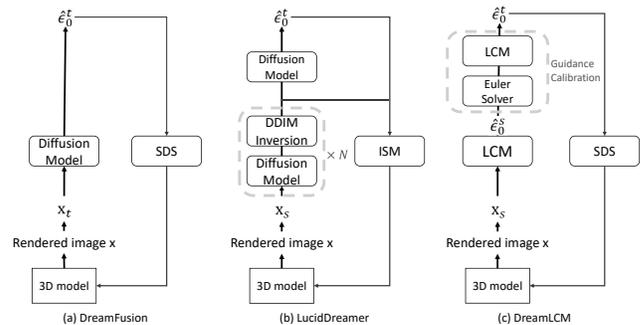
• Information systems → Multimedia content creation.

## KEYWORDS

Text-to-3D Generation, Diffusion Model, Gaussian Splatting, Latent Consistency Model

## 1 INTRODUCTION

Recent advancements in Diffusion Models (DMs) [34, 35, 37] have made progress in satisfying the needs of synthesizing high-quality images given text descriptions. Besides, by training on large-scale image datasets [38] where images are coupled with detailed texts, DMs achieve powerful ability in generating all kinds of 3D contents conditioned on the given text prompts. Existing works [3, 17, 22, 23, 27, 32, 44, 47] have been proposed to apply well-trained diffusion models to the task of text-to-3D generation. Conditioned on text prompts, DMs can generate guidance information in the latent



**Figure 1: Illustration of different guidance generation approaches.**  $x$  and  $\hat{\epsilon}$  indicates the rendered image and the guidance, respectively. (a) SDS generates guidance via a single diffusion model while producing over-smooth results. (b) LucidDreamer utilizes the DDIM inversion technique, forwarding Diffusion Models multiple times where  $N = \{2, 3, 4, 5\}$ . (c) The proposed DreamLCM method incorporates LCM as the guidance model. We also propose a guidance calibration strategy that uses Euler Solver to refine the guidance  $\hat{\epsilon}_0^s$  to  $\hat{\epsilon}_0^t$ . Our method generates higher-quality guidance compared to (a) and (b).

space. This latent guidance can be utilized to supervise the carving process of the target 3D objects. Thus, this alternative approach tackles the challenge of 3D object generation without large-scale 3D models for training. For example, the Score Distillation Sampling (SDS) objective is introduced in DreamFusion [32] to leverage the robust prior knowledge acquired by text-to-image diffusion models [35, 37]. The SDS backpropagates the gradients from the 2D diffusion model to 3D objects and bridges the gap between the diffusion models and the 3D representations, as shown in Fig. 1(a). The acquired prior knowledge is utilized to optimize the 3D objects represented by Neural Radiance Fields (NeRF) [28] conditioned on a single text prompt.

However, SDS is limited in generating fine details as it produces over-smooth results. This effect has been noted by previous works [17, 22, 45]. These works attempt to improve SDS and achieve good results in increasing the quality of 3D models. For instance, ProlificDreamer [45] optimizes multiple 3D models simultaneously. These models are mutually benefited and merged by finetuning a LoRA model [12]. The LoRA thus reserves the details of the 3D model. Nevertheless, extra resources are needed to regenerate the lost details. In this paper, we think that the problem of the over-smooth issue stems from two factors. Firstly, the guidance of SDS is derived by the DDPM [9] single-step inference, which leads to low-quality guidance and blurred details. Secondly, the rendered images of the target 3D object act as conditions and are fed into DMs

Unpublished working draft. Not for distribution.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted by ACM, provided that the copies are not made for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ACM MM, 2024, Melbourne, Australia

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

https://doi.org/10.1145/nmmmmmmmmmmmm

https://doi.org/10.1145/nmmmmmmmmmmmm

117 after adding random noises. Besides, DMs need to sample timesteps  
 118 randomly for the diffusion process. The randomness from both  
 119 the added noises and timesteps directly results in the *inconsistent*  
 120 *guidance* between different iterations. This inconsistency ultimately  
 121 averages the details of 3D models and leads to the over-smooth  
 122 issue.

123 This paper endeavors to tackle the over-smooth issue by han-  
 124 dling the factors above accordingly. We propose a novel approach to  
 125 generate clear guidance for the *low-quality guidance issue*. Inspired  
 126 by Latent Consistency Model (LCM) [26], we propose **DreamLCM**  
 127 by incorporating LCM as a guidance model to fully utilize the capa-  
 128 bility to generate high-quality guidance in a single-step inference.  
 129 Notably, LCM generates high-quality images in a single-step infer-  
 130 ence, rather than gradually approaching to the origin along the  
 131 probably flow ODE (PF-ODE) trajectory [42] by multi-step inference  
 132 like DDIM [40]. Therefore, DreamLCM merely predicts guidance  
 133 of a rendered image by directly denoising the noisy latent from an  
 134 arbitrary timestep along the PF-ODE trajectory to keep fine details  
 135 of the target 3D models. For the *inconsistency issue*, we observe  
 136 that LCM generates consistent guidance regardless of the random-  
 137 ness. To solve the issue, a similar method has been proposed in  
 138 LucidDreamer [22], which uses the DDIM inversion technique [40]  
 139 to improve the consistency of the guidance. However, different  
 140 from LucidDreamer, the proposed DreamLCM method provides  
 141 two merits: 1) DreamLCM only needs a single-step inference to  
 142 compute the guidance while LucidDreamer forwards the U-Net [36]  
 143 in DM multiple times; 2) DreamLCM keeps the original SDS loss.  
 144 Since LCM can resolve the two issues causing the over-smooth  
 145 issue, there is no need to change the loss forms. On the contrary,  
 146 LucidDreamer utilizes a complex objective function to adapt the  
 147 DDIM Inversion method. The difference is shown in Fig. 1(b)(c).

148 In addition, to further improve generation quality, we propose  
 149 two novel strategies, i.e., *Guidance Calibration* and *Dual Timestep*  
 150 *Strategy*. For *Guidance Calibration*, we propose a two-stage strat-  
 151 egy that repeats the perturbing and denoising steps to calibrate  
 152 the distribution of the guidance. In this way, the disturbing infor-  
 153 mation can be gradually removed. In the first stage, we perturb  
 154 a rendered image and predict the corresponding guidance. This  
 155 guidance is consistent with the rendered image, i.e., the 3D object,  
 156 as the added noise and timestep are small. In the second stage, we  
 157 run a discretization step of a numerical ODE solver, where we use  
 158 the Euler Solver to obtain a latent with relatively large noises. The  
 159 large noises and timestep can provide a more reasonable optimiza-  
 160 tion direction. Consequently, the calibrated guidance is ensured to  
 161 be consistent with both the rendered image and the highest data  
 162 density conditioned on the text prompt, effectively improving the  
 163 guidance’s quality. We calculate SDS loss and back-propagate the  
 164 gradient using the calibrated guidance to optimize the 3D models.  
 165 For *Dual Timestep Strategy*, we utilize the timestep sampling strat-  
 166 egy to enable dreamLCM to optimize the geometry and appearance  
 167 of 3D objects in separate phases. In particular, in the initial phase,  
 168 we apply large timesteps to guide the 3D model in producing large  
 169 deformations. In this case, DreamLCM tends to optimize geometry,  
 170 where the position of Gaussian Splatting is greatly updated. In the  
 171 refinement phase, we use small timesteps to optimize the appear-  
 172 ance because small timesteps help DreamLCM generate guidance  
 173 with fine details. Besides, we sample monotonically decreasing  
 174

timesteps to increase the consistency of the guidance. Overall, our  
 proposed dual timestep strategy is the combination of the timestep  
 strategy in HiFA [49] and ProlificDreamer [45].

We apply the Gaussian Splatting [18] as the 3D representation  
 to form the 3D target objects. The proposed DreamLCM achieves  
 the state-of-the-art results. As shown in Fig. 3, we can see that  
 DreamLCM generates high-quality 3D objects with fine details.  
 Besides, our model trains end-to-end, reducing training costs and  
 maintaining a streamlined training pipeline. Overall, our contribu-  
 tions can be summarized as follows:

- We resolve the over-smooth issue of SDS in a new perspective  
 by proposing DreamLCM. We analyze the two weaknesses  
 in the generated guidance of diffusion models, i.e., low qual-  
 ity and low consistency. In response to the two issues, we  
 incorporate LCM as our guidance model to make full use  
 of the ability in LCM and generate high-quality, consistent  
 guidance in a single inference step.
- We propose two novel strategies to further improve the qual-  
 ity of the guidance for 3D generation. A guidance calibra-  
 tion strategy is proposed, using Euler solver to obtain an  
 improved sample, which subsequently generates calibrated  
 guidance to help 3D models converge accurately. Besides, a  
 dual timestep strategy is proposed, enabling DreamLCM to  
 optimize the geometry and the appearance in two phases.  
 We prove the effectiveness of the two strategies in Sec.6.4.
- We conduct experiments to demonstrate that DreamLCM  
 significantly outperforms the state-of-the-art methods in  
 terms of both quality and training efficiency.

## 2 RELATED WORK

### 2.1 Diffusion Models

Diffusion Models(DMs) have emerged as powerful tools for image  
 generation [10, 29, 31, 34, 40, 42], excelling in denoising noise-  
 corrupted data and estimating data distribution scores. The stable  
 ability of DMs for generating high-quality images led to multiple  
 applications in various domains, including video [8, 11, 19] and  
 3D [32, 44], *etc.* During inference, these models employ reverse dif-  
 fusion processes to gradually denoise data points and generate sam-  
 ples. In comparison to Variational Autoencoders (VAEs) [20, 39] and  
 Generative Adversarial Networks (GANs) [7], diffusion models of-  
 fer enhanced training stability and likelihood estimation. However,  
 their sampling efficiency is often hindered. Discretizing reverse-  
 time SDE [6, 42] or ODE [42] are proposed to handle the challenge,  
 various techniques such as ODE solvers [24, 25, 40, 48], adaptive  
 step size solvers [14], and predictor-corrector methods [42] have  
 been proposed. Notably, the Latent Diffusion Model(LDM) [35]  
 conducts forward and reverse diffusion processes in the latent  
 data space, leading to more efficient computation. The Consistency  
 model [26, 41] demonstrates promising potential as a rapid sam-  
 pling generative model for generating high-quality images in a  
 single-step inference. In this paper, we transfer the ability of LCM  
 to text-to-3D task to generate high-quality guidance. Meanwhile,  
 we use Euler Solver as the numerical ODE Solver to further calibrate  
 the guidance for higher quality.

## 2.2 Text-to-3D Generation.

This task targets generating 3D contents from given text prompts. The 3D contents are parameterized by various 3D representations, including implicit representations [1, 2, 5, 16], 3D Gaussians [4, 21, 22, 43, 46], *etc.* Existing methods includes 3D generative methods [15, 30]. However, these methods can only generate objects within limited categories due to the lack of large-scale datasets. Our method uses DMs to guide the 3D optimization. DreamField [13] represent pioneering efforts in training Neural Radiance Fields (NeRF) [28] with guidance from CLIP [33]. Dreamfusion [32] firstly employs Score Distillation Sampling (SDS) to distil 3D assets from pretrained text-to-image diffusion models. SDS has become integral to subsequent works, with endeavours aiming at enhancing 3D representations, addressing inherent challenges such as the Janus problem, and mitigating the over-smooth effect observed in SDS. Recent studies like ProlificDreamer [45], HiFA [49], and LucidDreamer [22] have made significant strides in improving the SDS loss. Concurrent methods such as CSD [47] and NFSD [17] provide empirical solutions to enhance SDS. In our novel approach, DreamLCM, we resolve the over-smooth problem in a new perspective of the guidance model, showing that it is possible to generate high-quality 3D models without any alterations to SDS.

## 3 REVISITING CONSISTENCY MODELS

The core idea of the Consistency Model (CM) and Latent Consistency Model (LCM) is to learn a function that maps any points on a trajectory of PF-ODE [42] to that trajectory origin, i.e., the solution of PF-ODE. The trajectory origin indicates the real data distribution region, which has the highest data density. Besides, LCM extends the denoising process to the latent space. LCM predicts the solution of PF-ODE by introducing a consistency function in a single-step inference. LCM is a text-to-image DM  $f_\phi$  parameterized by  $\phi$ . The objective is to fulfill the mapping:  $f_\phi(\mathbf{x}_t, \mathbf{y}, t) \rightarrow \mathbf{x}_0$ , where  $\mathbf{x}_t$  is the noisy latent while  $\mathbf{y}$  is the text prompt. The self-consistency property of LCM is expressed in Eq. (1) as

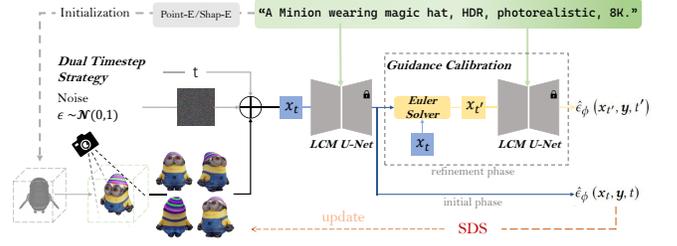
$$f_\phi(\mathbf{x}_t, \mathbf{y}, t) = f_\phi(\mathbf{x}_{t'}, \mathbf{y}, t'), \forall t, t' \in [\delta, T], \quad (1)$$

where  $\delta$  is a fixed small positive number. The formula shows the consistency of the perturbed images between different timesteps.

Overall, LCM has two benefits: 1) generating  $f_\phi(\mathbf{x}_t, \mathbf{y}, t)$  with high quality in a single-step inference. 2) different  $\mathbf{x}_t$  between different  $t$  generate consistent guidance. We attribute the over-smooth issue in SDS loss to two factors in Sec. 1. The first is the *low-quality guidance* issue, which can be resolved by utilizing LCM to generate high-quality guidance. The second is the *inconsistency issue*. The issue is mitigated because the guidance generated via LCM is consistent between different timesteps. Therefore, we incorporate LCM into the text-to-3D task as the guidance model. The guidance generated by LCM exhibits high quality and high consistency.

## 4 METHOD

In this section, we present DreamLCM for high-quality text-to-3D synthesis. First, we formulate the entire 3D generation process and analyze the issues in recent works. Then, we propose DreamLCM, Guidance Calibration, and Dual Timestep Strategy. We explain how these methods benefit the generation quality.



**Figure 2: Illustration of DreamLCM.** DreamLCM initializes the 3D model  $\theta$  via text-to-3D generator [15, 30]. We utilize the proposed timestep strategy to divide the training into two phases. In the initial phase, we directly generate guidance via a single LCM network. In the refinement phase, we utilize another LCM network and an Euler Solver to calibrate the guidance. We calculate the original SDS loss to update  $\theta$ .

## 4.1 The Problem Definition

Dreamfusion [32] proposes a general framework for the text-to-3D generation task. It has two important components. The first is a 3D representation, e.g., NeRF [28], 3D Gaussian Splatting [18], that is parameterized by  $\theta$  for depicting the target 3D object  $\Theta$ . The second is a pretrained text-to-image diffusion model for providing guidance information and supervising the target  $\Theta$ . To bridge the 3D object and its guidance model, a differentiable renderer  $g$  is utilized to obtain the rendered image  $\mathbf{x}$ , which is formulated as  $\mathbf{x} = g(\theta, c)$  with camera pose  $c$ . Then  $\mathbf{x}$  is fed into a VAE encoder [20] and perturbed by noise  $\epsilon$ . Here, we denote the latent embedding as  $\mathbf{x}$  for simplification. Given noisy latent  $\mathbf{x}_t$ , timestep  $t$ , and text prompt  $\mathbf{y}$  as inputs, the guidance model predicts guidance gradients for updating the 3D object. The guidance prediction of DM can be expressed by two equivalent forms, i.e.,  $\epsilon$ -prediction  $\hat{\epsilon}_\phi(\mathbf{x}_t, \mathbf{y}, t)$  and  $\mathbf{x}$ -prediction  $\hat{\mathbf{x}}_0^t$  in Eq (2) following SDS [32].

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\text{SDS}}(\phi, g) &= \mathbb{E}_{t, \epsilon, c} \left[ \omega(t) \left( \hat{\epsilon}_\phi(\mathbf{x}_t, \mathbf{y}, t) - \epsilon \right) \frac{\partial g(\theta, c)}{\partial \theta} \right], \\ &= \mathbb{E}_{t, \epsilon, c} \left[ \frac{\omega(t)}{\gamma(t)} (\mathbf{x}_0 - \hat{\mathbf{x}}_0^t) \frac{\partial g(\theta, c)}{\partial \theta} \right]. \end{aligned} \quad (2)$$

Here  $\hat{\mathbf{x}}_0^t = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}_\phi(\mathbf{x}_t, \mathbf{y}, t)}{\sqrt{\bar{\alpha}_t}}$ ,  $\sqrt{\bar{\alpha}_t}$  is the noise weight, which shows that  $\hat{\epsilon}_\phi(\mathbf{x}, \mathbf{y}, t)$  and  $\hat{\mathbf{x}}_0^t$  are equivalent, and we consider them both as guidance.  $\omega(t)$  is a weighting function that depends on the timestep  $t$ .  $\gamma(t) = \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}}$ . The gradient leads the 3D model closer to the corresponding text prompt.

Previous works ProlificDreamer [45] and LucidDreamer [22] observe that SDS generates over-smooth 3D models. We attribute this issue to two factors: 1) DMs [35, 37] generate low quality guidance because  $\hat{\mathbf{x}}_0^t$  are obtained from DDPM single-step inference [9], as shown in Fig. 1(a); 2) DMs are sensitive to the randomness in noise  $\epsilon$  and timestep  $t$ . Especially, a large  $t$  is hard to generate  $\hat{\mathbf{x}}_0^t$  consistent with  $\mathbf{x}_0$ , which averages the appearance of the 3D models during optimization. Overall, the weakness in DMs generates poor guidance, resulting in over-smooth outcomes.

## 4.2 DreamLCM

The proposed DreamLCM method aims at resolving the aforementioned over-smooth issue by incorporating LCM and further enhancing the generation quality by proposing two effective strategies, i.e., Guidance Calibration and a dual timestep strategy. The overall framework is shown in Fig. 2. The entire DreamLCM approach is depicted in Algorithm 1.

For the *low-quality guidance issue*, DreamLCM utilizes the powerful ability of LCM to generate high-quality guidance. LCM trains a function  $f_\phi$  in Eq. 1, which can be seen as guidance, to map any  $\mathbf{x}_t$  to its PF-ODE trajectory origin, i.e.. Consequently, LCM is capable of generating high-quality guidance in a single-step inference.

For the *inconsistency issue*, we utilize the important property of LCM in Eq. 1, highlighting the consistency of guidance between different timesteps. When guided by DMs, rendered image  $\mathbf{x}_0$  is added random noise  $\epsilon$ , which is further weighted by different timesteps  $t$ . The randomness in  $\epsilon$  and  $t$  directly results in inconsistent  $\hat{\mathbf{x}}_0^t$ , eventually resulting in a *feature-average* outcome. However, due to LCM's property, LCM can generate consistent  $\hat{\mathbf{x}}_0^t$  regardless of the randomness.

Overall, given a timestep  $s$ , we integrate LCM and SDS by generating guidance  $\hat{e}_\phi(\mathbf{x}_s, \mathbf{y}, s)$  via LCM to calculate the the SDS loss :

$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\phi, \mathbf{g}) = \mathbb{E}_{s, \epsilon, c} \left[ w(s) \left( \hat{e}_\phi(\mathbf{x}_s, \mathbf{y}, s) - \epsilon \right) \frac{\partial \mathbf{g}(\theta, c)}{\partial \theta} \right], \quad (3)$$

where  $\hat{e}_\phi(\mathbf{x}_s, \mathbf{y}, s)$  is obtained in a single-step inference with high quality and fine details. It enables the 3D model to have fine details, mitigate the over-smooth issue, and save training costs.

Unfortunately, it is difficult to resolve the two issues perfectly due to the nature of diffusion models. The images generated by LCM with a single-step inference are always blurred, and the high-quality images are derived from four-step inferences iteratively. This fact is also stated in LCM [26]. For the first issue, if we directly utilize LCM's single-step inference results, the guidance would be blurred and not conducive to generating high-quality 3D models. Therefore, we further resolve this weakness by proposing a *guidance calibration* strategy. For the second issue, during the training of LCM, the two noisy latents in Eq. (1) are limited to be on the same PF-ODE trajectory, rather than two arbitrary noisy latents. We follow this protocol during the inference by fixing the noise  $\epsilon'$  to perturb the rendered image, reducing the randomness in noise. Besides, we propose a decreasing timestep strategy, where we utilize monotonically decreasing timesteps during training, to reduce the randomness in timesteps.

**Guidance Calibration.** We further dive into the principle when LCM generates guidance. We first review that in DMs, the denoising process follows a reverse stochastic differential equation(SDE):

$$d\mathbf{x} = -\dot{\sigma}_t \sigma_t \nabla \log p_t(\mathbf{x}) dt + \sqrt{\dot{\sigma}_t \sigma_t} d\mathbf{w}, \quad (4)$$

where  $p_t(\mathbf{x}_t) \sim \mathcal{N}(\mathbf{x}_0, \sigma_t^2 \mathbf{I})$ ,  $\sigma_t$  varies along timestep  $t$ ,  $\dot{\sigma}_t$  is the time derivative of  $\sigma_t$ ,  $\mathbf{w}$  is the standard Wiener process and  $\nabla \log p_t(\mathbf{x})$  is the score function which indicates the direction towards highest data density. And there exists a corresponding reverse ordinary deterministic equation(ODE) defined below:

$$d\mathbf{x} = -\dot{\sigma}_t \sigma_t \nabla \log p_t(\mathbf{x}) dt. \quad (5)$$

where  $\nabla \log p_t(\mathbf{x})$  is estimated as  $-\frac{1}{\sqrt{1-\alpha_t}} \hat{e}_0^t$ . LCM can map any  $\mathbf{x}_t$  on a trajectory of the PF-ODE to the origin  $\mathbf{x}_0^*$ , which indicates the highest-data-density region. However, the mapped origin, i.e.,  $\hat{\mathbf{x}}_0^t$  derived from single-step inference is always shifted. We attribute the shifting to the insufficient training of  $f_\phi$ . Our goal is to calibrate  $\hat{\mathbf{x}}_0^t$  to get closer to  $\mathbf{x}_0^*$ . We consider the insufficient training in LCM and rationally assume that the denoising process of LCM follows a smooth PF-ODE trajectory with a small slope. This assumption allows us to calibrate the guidance from the perspective of PF-ODE.

Based on the analysis, we propose our guidance calibration strategy, which is a two-stage strategy. We repeat the perturbing and denoising steps to calibrate the guidance distribution. In the first stage, given a perturbed sample  $\mathbf{x}_s$  at timestep  $s$ , we first predict guidance  $\hat{e}_\phi(\mathbf{x}_s, \mathbf{y}, s)$ , where  $s$  is set to a small number to make the guidance more consistent with  $\mathbf{x}_0$  than large  $s$ . In the second stage, since the denoising process of LCM follows PF-ODE, we run a discretization step of a numerical ODE solver. we use Euler Solver to get another sample  $\mathbf{x}_t$ :

$$\mathbf{x}_t = \frac{\sqrt{1-\sigma_s^2} \mathbf{x}_s + (\sigma_t - \sigma_s) \hat{e}_\phi(\mathbf{x}_s, \mathbf{y}, s)}{\sqrt{1-\sigma_t^2}} \quad (6)$$

where  $t > s$ . We then fed  $\mathbf{x}_t$  to LCM network to obtain the final calibrated guidance  $\hat{e}_\phi(\mathbf{x}_t, \mathbf{y}, t)$ . Compared to  $\hat{e}_\phi(\mathbf{x}_s, \mathbf{y}, s)$ , the guidance  $\hat{e}_\phi(\mathbf{x}_t, \mathbf{y}, t)$  has two advantages: 1) it is consistent with the original rendered image  $\mathbf{x}_0$  because the Euler solver makes  $\mathbf{x}_t$  and  $\mathbf{x}_s$  on the same PF-ODE trajectory; 2) large timestep  $t$  provides more reasonable optimization direction conditioned on  $\mathbf{y}$ , leading  $\hat{\mathbf{x}}_0^t$  closer to  $\mathbf{x}_0^*$ . Overall,  $\hat{\mathbf{x}}_0^t$  is ensured to be consistent with both the rendered image and the PF-ODE trajectory origin conditioned on the text prompt, effectively improving the quality of the guidance. We optimize  $\theta$  using the final guidance  $\hat{\mathbf{x}}_0^t$  to calculate SDS loss following Eq. (7) as

$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\phi, \mathbf{g}) = \mathbb{E}_{t, \epsilon, c} \left[ w(t) \left( \hat{e}_\phi(\mathbf{x}_t, \mathbf{y}, t) - \epsilon \right) \frac{\partial \mathbf{g}(\theta, c)}{\partial \theta} \right]. \quad (7)$$

**Dual Timestep Strategy.** In this paper, we incorporates 3D Gaussians [18] as the 3D representation, which requires initialization models, i.e., PointE [30], Shap-E [15] to initialize the geometry. However, these models sometimes initialize badly, especially when given complex text prompts. Thus, properly updating the geometry positions of the 3D model is crucial in 3D generation. We propose a two-phase strategy, optimizing the geometry and appearance of 3D objects in separate phases. In the initial phase, we use large timesteps to predict large deformations to the 3D model since large timesteps keep less information in the rendered image. As a result, the guidance includes global geometry features, leading DreamLCM to optimize the geometry, where the position of Gaussian Splatting is greatly updated. In the refinement phase, we use small timesteps to optimize the appearance because small timesteps keep more information on the rendered image, generating guidance with fine local features.

We propose a dual timestep strategy combining the decreasing timestep strategy with the two-phase strategy. Specifically, we define a cut-off iteration  $T_{cut}$  and a cut-off timestep  $t_{cut}$ , in each

**Algorithm 1** DreamLCM

---

```

1: Initialization: 3D model parameters  $\theta$ , training iteration  $n$ ,
   LCM network  $\phi$  denoising timestep from  $N_{min}$  to  $N_{max}$ , cut-off
   iteration  $T_{cut}$  and timestep  $t_{cut}$ , text prompt  $\mathbf{y}$ , fixed noise  $\epsilon'$ .
2: for  $i = [0, \dots, n - 1]$  do
3:   if  $i \leq T_{cut}$  then
4:      $t_{max} \leftarrow N_{max}, t_{min} \leftarrow t_{cut}, t_{interval} \leftarrow T_{cut}, id \leftarrow i$ 
5:   else
6:      $t_{max} \leftarrow t_{cut}, t_{min} \leftarrow N_{min}, t_{interval} \leftarrow n - T_{cut},$ 
        $id \leftarrow i - t_{cut}$ 
7:   end if
8:   Sample: camera pose  $c, \mathbf{x}_0 = g(\theta, c)$ 
9:    $s \leftarrow t_{max} - (t_{max} - t_{min}) \sqrt{id/t_{interval}}, t \leftarrow 2s$ 
10:   $\mathbf{x}_s \leftarrow \mathbf{x}_0 + \sigma_s \epsilon'$ 
11:  predict  $\hat{\epsilon}_\phi(\mathbf{x}_s, \mathbf{y}, s)$ 
12:  if  $i \leq T_{cut}$  then
13:    calculate SDS loss:
14:     $\nabla_{\theta} L_{SDS} = \omega(s) \left( \hat{\epsilon}_\phi(\mathbf{x}_s, \mathbf{y}, s) - \epsilon \right)$ , update  $\theta$ .
15:  else
16:    use Euler Solver to obtain  $\mathbf{x}_t$ .
17:    predict  $\hat{\epsilon}_\phi(\mathbf{x}_t, \mathbf{y}, t)$  then calculate SDS loss:
18:     $\nabla_{\theta} L_{SDS} = \omega(t) \left( \hat{\epsilon}_\phi(\mathbf{x}_t, \mathbf{y}, t) - \epsilon \right)$ , update  $\theta$ .
19:  end if
20: end for

```

---

stage, the timestep is calculated as follows :

$$t = t_{max} - (t_{max} - t_{min}) \sqrt{T/N}, \quad (8)$$

where  $T$  and  $N$  are the current iteration and total iteration. For the first  $T_{cut}$  iterations, we optimize geometry using timesteps larger than  $t_{cut}$ . For the remaining iterations, we use timesteps less than  $T_{cut}$  to optimize appearance. Overall, we can see that the timestep strategy in HiFA [49] and ProlificDreamer [45] are two special cases of our timestep strategy. Experiments demonstrate that the strategy can generate high-quality 3D models with fine geometry and appearance, as shown in Fig. 5.

## 5 DISCUSSION

Similar to the proposed DreamLCM method, ProlificDreamer [45] and LucidDreamer [22] targets resolving the over-smooth issue in SDS. They refine the SDS loss with different loss functions to alleviate the over-smoothed and over-saturated results based on SDS. We will revisit these two losses to show the relationship between DreamLCM and the two works and demonstrate that our work is more effective than theirs.

**ProlificDreamer** is based on the SDS loss. It handles the over-smooth issue by training an additional LoRA [12] network denoted as  $\epsilon_{LoRA}$ . ProlificDreamer optimizes multiple 3D models simultaneously. It aggregates and estimates their distributions by finetuning  $\epsilon_{LoRA}$ . The VSD loss for the  $i_{th}$  3D model is as follows:

$$\nabla_{\theta^{(i)}} \mathcal{L}_{VSD}(\theta^{(i)}) = \mathbb{E}_{t, \epsilon, c} \left[ \omega(t) \left( \hat{\epsilon}_\phi(\mathbf{x}_t^{(i)}, \mathbf{y}, t) - \hat{\epsilon}_{LoRA}(\mathbf{x}_t^{(i)}, \mathbf{y}, t, c) \right) \frac{\partial g(\theta^{(i)}, c)}{\partial \theta^{(i)}} \right], \quad (9)$$

where  $\mathbf{x}_t^{(i)}$  is the rendered image of the  $i_{th}$  3D model and  $c$  is the camera condition.  $\hat{\epsilon}_{LoRA}(\mathbf{x}_t^{(i)}, \mathbf{y}, t, c)$  indicates the distribution of the rendered image. When ProlificDreamer optimizes one 3D model, the distribution of the rendered image can be estimated as  $\epsilon$ , where  $\epsilon$  is the noise added to the rendered image. We consider the SDS gradient as the vector starting from  $\epsilon$  and  $\hat{\epsilon}_{LoRA}$  to  $\hat{\epsilon}_\phi$ . Since  $\hat{\epsilon}_{LoRA}$  contains information from multiple 3D models,  $\hat{\epsilon}_{LoRA}$  is a steadier and more robust starting point than  $\epsilon$ , averaging the random and inconsistent features in the optimization process of each 3D model.

We observe that the essential problem is the randomness and inconsistency when optimizing a single 3D model. Besides,  $\hat{\epsilon}_{LoRA}$  introduces extra parameters and trains several 3D models simultaneously, leading to high training costs. However, the proposed DreamLCM method incorporates LCM as the guidance model, greatly mitigating the inconsistent issue when optimizing one 3D model. As a result, there is no need for DreamLCM to train another  $\hat{\epsilon}_{LoRA}$  to decrease the training costs for generating high-quality 3D models.

**LucidDreamer** proposes ISM [22] loss, which employs DDIM Inversion to enhance the quality and consistency of the guidance. Specifically, it predicts a invertible noisy latent trajectory  $\{\mathbf{x}_{\delta_T}, \mathbf{x}_{2\delta_T}, \dots, \mathbf{x}_t\}$ , iteratively following Eq. (10),

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \hat{\mathbf{x}}_0^s + \sqrt{1 - \bar{\alpha}_t} \epsilon_\phi(\mathbf{x}_s, \theta, s), \quad (10)$$

where  $s = t - \delta t$ . The guidance is obtained by a multi-step DDIM denoising process i.e., iterating

$$\tilde{\mathbf{x}}_{t-\delta_T} = \sqrt{\bar{\alpha}_{t-\delta_T}} \left( \hat{\mathbf{x}}_0^t + \gamma(t - \delta_T) \epsilon_\phi(\mathbf{x}_t, \mathbf{y}, t) \right), \quad (11)$$

where  $\eta(t) = \frac{1 - \sqrt{\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}}$ . Next, by replacing  $\hat{\mathbf{x}}_0^t$  in Eq. (2) with  $\tilde{\mathbf{x}}_0$ , the SDS loss can be rewrote as  $\nabla_{\theta} \mathcal{L}(\theta) = \mathbb{E}_c \left[ \frac{\omega(t)}{\gamma(t)} \left( \mathbf{x}_0 - \tilde{\mathbf{x}}_0^t \right) \frac{\partial g(\theta, c)}{\partial \theta} \right]$ . LucidDreamer then unifies the iterative process in Eq. (10) and Eq. (11), proposing ISM loss as follows:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}(\theta) &= \mathbb{E}_{t, c} \left[ \frac{\omega(t)}{\gamma(t)} \left( \gamma(t) [\epsilon_\phi(\mathbf{x}_t, \mathbf{y}, t) - \epsilon_\phi(\mathbf{x}_s, \theta, s)] + \eta_t \right) \frac{\partial g(\theta, c)}{\partial \theta} \right] \\ &\approx \mathbb{E}_{t, c} \left[ \omega(t) \left( \epsilon_\phi(\mathbf{x}_t, \mathbf{y}, t) - \epsilon_\phi(\mathbf{x}_s, \theta, s) \right) \frac{\partial g(\theta, c)}{\partial \theta} \right]. \end{aligned} \quad (12)$$

where  $\eta_t$  includes a series of neighboring interval scores with opposing scales, which can be disregarded. ISM essentially substitutes DDPM single-step inference [9] for DDIM multi-step inference [40] to generate high-quality and high-fidelity guidance  $\tilde{\mathbf{x}}_0^t$ .

However, the multi-step inference needs to forward the U-Net [36] in DMs multiple times, increasing training costs. Moreover, we can see that a key improvement in ISM is the quality and consistency of the guidance. Compared to LucidDreamer, DreamLCM is capable of generating high-quality and high-consistency guidance in a single-step inference. Consequently, DreamLCM is more effective with fewer training costs.

To sum up, we observe that the principal cause of the over-smooth issue in SDS is the inadequate quality of the guidance. These two methods tackle the issue by utilizing extra resources, e.g., training multiple NeRFs and DDIM Inversions, which is time-consuming. Differently, DreamLCM can resolve the over-smooth issue by taking full advantage of LCM, while saving training costs.

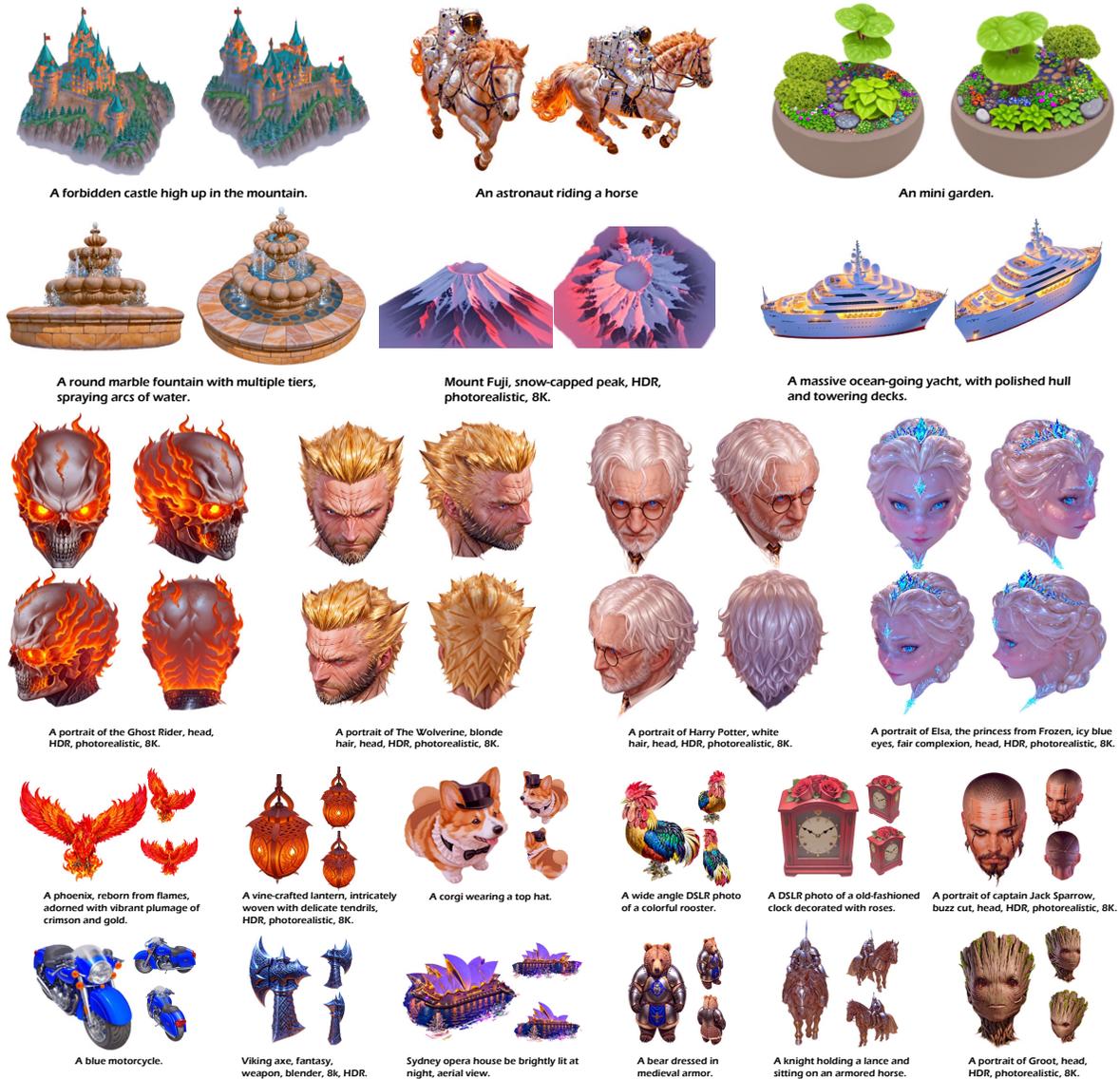


Figure 3: Examples generated by DreamLCM. We incorporate the Latent Consistency Model (LCM) as a guidance model, with two proposed strategies to further enhance the generation quality (See section 4 for details). DreamLCM generates high-quality results with fine details.

## 6 EXPERIMENTS

### 6.1 Implementation Details

We train our end-to-end network for 5000 iterations overall. We employ 3D Gaussian Splatting [18] as our 3D representation and 3D point cloud generation models Shap-E [15] and Point-E [30] for initialization. The rendering resolution is  $512 \times 512$ . As for the guidance calibration strategy, we use it in appearance optimization. We practically consider  $s = 350$  as the cut-off timestep. Unless stated otherwise, we train the first 1000 iters for geometry optimization using timesteps  $s$  fulfilling  $350 \leq s \leq 980$  and the remaining 4000 iters for appearance optimization using timesteps  $s$  fulfilling

$20 \leq s \leq 350$ . Since we assume that LCM follows a smooth PF-ODE, the interval between  $s$  and  $t$  is less limited. Practically, we choose  $t = 2s$ . We use SDS with a normal CFG scale of 7.5. All experiments are performed and measured with an RTX 3090 (24G) GPU. We train about 50 min per sample.

### 6.2 Text-to-3D Generation.

In Fig. 4, we show the generated results of DreamLCM. We generate all examples using the original LCM without LoRA and any finetuned checkpoints. We can see that DreamLCM can generate photo-realistic 3D objects with fine details. The 3D objects are creative and highly consistent with the text prompts. Especially, we

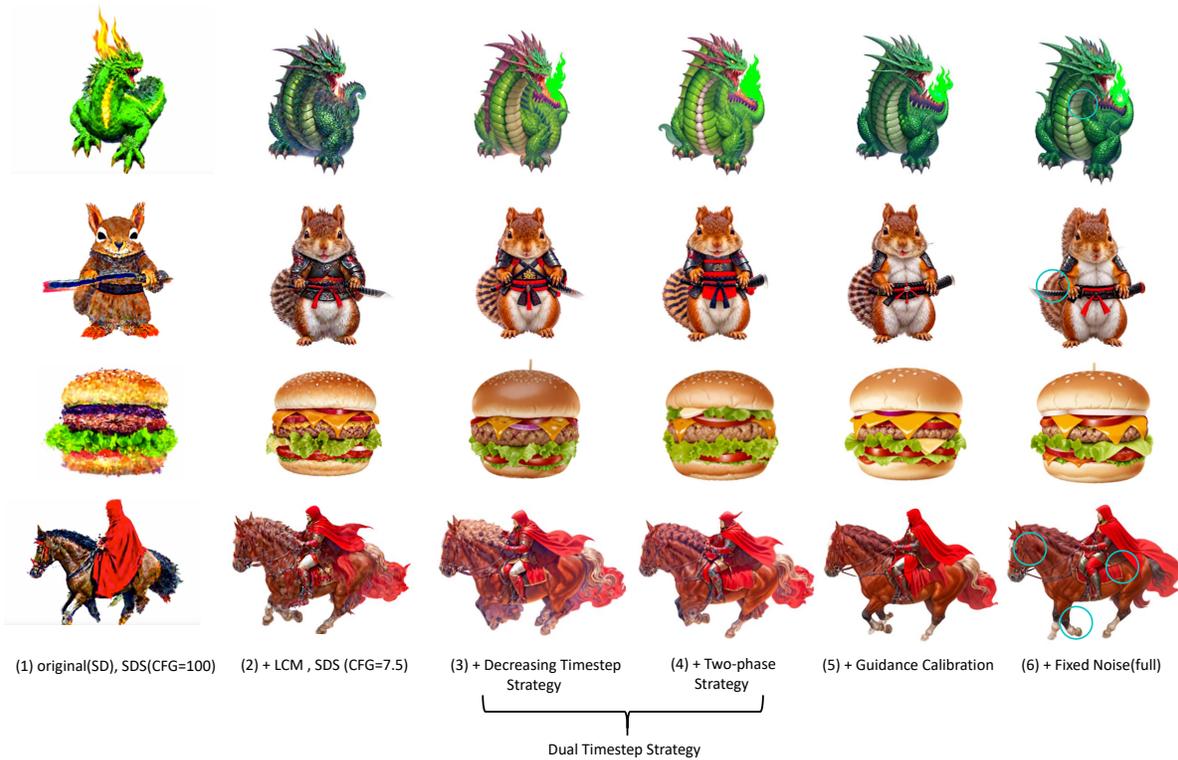


**Figure 4: Comparison with the state-of-the-art text-to-3D generation methods with Gaussian Splatting as 3D representations. Experiments show that the proposed DreamLCM generates photo-realistic 3D objects with high quality and fine details. The models generated by DreamLCM are more consistent with the text prompt. The training time is measured with a single RTX 3090 GPU.**

can see that DreamLCM can generate different and amazing avatar heads conditioned on text prompts, such as "A portrait of Harry Potter, white hair, head, HDR, photorealistic, 8K". Besides, the proposed DreamLCM method is good at generating objects conditioned on complex text prompts, like "the fuji mountain", "the massive yacht", and "the fountain". These examples demonstrate that DreamLCM well resolves the over-smooth issue in SDS. Besides, these examples show great potential in generating all kinds of complex objects with different LCM finetuned checkpoints.

### 6.3 Qualitative Comparison

We compare our method with the current SoTA baselines which generate 3D Gaussian Objects [22, 43, 46]. As shown in Fig. 4, our model generates more photo-realistic results than other works, exhibiting high quality and fine details. For example, "A portrait of a unicorn" is more photo-realistic, and the fur is more silky than the results from the other three approaches. As for the results conditioned on the text "A Spanish galleon", our model generates the most intact galleon, and the details of the hull are the finest



**Figure 5: Ablation Study of DreamLCM.** The proposed components are effective and can improve the text-to-3D generation quality. (1) The results of SDS with a large CFG scale of 100. (2) We incorporate LCM as a guidance model with a small CFG of 7.5. (3)(4) The results after adding the Dual Timestep Strategy. It includes two parts, the Decreasing Timestep Strategy to reduce the randomness in timesteps and the Two-phases Strategy to improve geometry. Both parts are effective. (5) The results after adding Guidance Calibration to further improve the generation quality. (6) We use fixed noise to perturb the samples to reduce the randomness in noises to improve the details. We highlight some improved details in cyan. The prompts corresponding to the four examples are "a green dragon breathing fire", "a squirrel in samurai armor wielding a katana", "a delicious hamburger" and "A warrior with red cape riding a horse".

among all the shown methods. Notably, compared to LucidDreamer, we generate higher quality objects with less training costs.

## 6.4 Ablation Study

Fig. 5 depicts the ablation experiments of different baseline methods. In Fig. 5(b) and (c), we utilize timesteps between 20 and 500 to generate high-quality images. Other settings are the same as the final settings 6.1. Starting from the original SDS loss, guided by Stable Diffusion [35], with a large CFG scale(100). We first incorporate LCM as the guidance model to demonstrate that LCM is a superior guidance model to DMs [35]. We can see that LCM makes a huge improvement in generation quality. We then add our Dual Timestep Strategy. We divide the strategy into two parts. We demonstrate the effectiveness of Decreasing Time Strategy and the Two-phase Strategy, as shown in Fig. 5(c) and (d). We can see that the hamburger adding Decreasing Time strategy shows a quality improvement. Based on (c), the hamburger adding the two-phase strategy shows a geometric advancement of the bread at the bottom of the hamburger. Besides, due to the two-phase strategy, the warrior example is deformed to be less close to the horse, since

this 3D model is initialized by the prompt "a wolf". Then, we add the guidance calibration strategy which smooths the appearance and improves the details, making the objects more photo-realistic. Finally, we add fixed noises to improve the consistency of guidance between different timesteps. As shown in Fig. 5, the ability to improve details is demonstrated in cyan, such as the eye, legs, and cushion on the horse back in the warrior sample, the katana in the squirrel sample and the shadow in the dragon sample.

## 7 CONCLUSION

In this paper, we propose DreamLCM method to improve the text-to-3D object task. We incorporate LCM as a guidance model to generate high-quality guidance to resolve the two factors that cause the over-smooth issue. Besides, we introduce two techniques, i.e., Guidance Calibration and Dual Timestep Strategy, to further improve the generation quality. Experiments show superior performance of our method. Our method achieves state-of-the-art results in both generation and training efficiency.

## REFERENCES

- [1] Titas Anciukevičius, Zexiang Xu, Matthew Fisher, Paul Henderson, Hakan Bilen, Niloy J. Mitra, and Paul Guerrero. 2023. RenderDiffusion: Image Diffusion for 3D Reconstruction, Inpainting and Generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 12608–12618.
- [2] Hansheng Chen, Jiatao Gu, Anpei Chen, Wei Tian, Zhuowen Tu, Lingjie Liu, and Hao Su. 2023. Single-Stage Diffusion NeRF: A Unified Approach to 3D Generation and Reconstruction. arXiv:2304.06714
- [3] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. 2023. Fantasia3D: Disentangling Geometry and Appearance for High-quality Text-to-3D Content Creation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 22246–22256.
- [4] Zilong Chen, Feng Wang, Yikai Wang, and Huaping Liu. 2024. Text-to-3D using Gaussian Splatting. arXiv:2309.16585
- [5] Yen-Chi Cheng, Hsin-Ying Lee, Sergey Tulyakov, Alexander Schwing, and Liangyan Gui. 2023. SDFusion: Multimodal 3D Shape Completion, Reconstruction, and Generation. arXiv:2212.04493
- [6] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. 2022. Score-Based Generative Modeling with Critically-Damped Langevin Diffusion. arXiv:2112.07068
- [7] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Networks. arXiv:1406.2661
- [8] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. 2022. Imagen Video: High Definition Video Generation with Diffusion Models. arXiv:2210.02303
- [9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. 6840–6851.
- [10] Jonathan Ho and Tim Salimans. 2022. Classifier-Free Diffusion Guidance. arXiv:2207.12598
- [11] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet. 2022. Video Diffusion Models. arXiv:2204.03458
- [12] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021).
- [13] Ajay Jain, Ben Mildenhall, Jonathan T. Barron, Pieter Abbeel, and Ben Poole. 2021. Zero-Shot Text-Guided Object Generation with Dream Fields. arXiv (2021).
- [14] Alexia Jolicœur-Martineau, Ke Li, Remi Piche-Taillefer, Tal Kachman, and Ioannis Mitliagkas. 2021. Gotta Go Fast When Generating Data with Score-Based Models. *ArXiv abs/2105.14080* (2021).
- [15] Heewoo Jun and Alex Nichol. 2023. Shap-E: Generating Conditional 3D Implicit Functions. arXiv:2305.02463
- [16] Animesh Karnewar, Niloy J. Mitra, Andrea Vedaldi, and David Novotny. 2023. HoloFusion: Towards Photo-realistic 3D Generative Modeling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 22976–22985.
- [17] Oren Katzir, Or Patashnik, Daniel Cohen-Or, and Dani Lischinski. 2023. Noise-Free Score Distillation. arXiv:2310.17590
- [18] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. arXiv:2308.04079
- [19] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. 2023. Text2Video-Zero: Text-to-Image Diffusion Models are Zero-Shot Video Generators. arXiv:2303.13439
- [20] Diederik P Kingma and Max Welling. 2022. Auto-Encoding Variational Bayes. arXiv:1312.6114
- [21] Xinhai Li, Huaibin Wang, and Kuo-Kun Tseng. 2023. GaussianDiffusion: 3D Gaussian Splatting for Denoising Diffusion Probabilistic Models with Structured Noise. arXiv:2311.11221
- [22] Yixun Liang, Xin Yang, Jiantao Lin, Haodong Li, Xiaogang Xu, and Yingcong Chen. 2023. LucidDreamer: Towards High-Fidelity Text-to-3D Generation via Interval Score Matching. arXiv:2311.11284
- [23] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaoohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. 2023. Magic3D: High-Resolution Text-to-3D Content Creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 300–309.
- [24] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. 2022. DPM-Solver: A Fast ODE Solver for Diffusion Probabilistic Model Sampling in Around 10 Steps. *ArXiv abs/2206.00927* (2022).
- [25] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. 2023. DPM-Solver++: Fast Solver for Guided Sampling of Diffusion Probabilistic Models. arXiv:2211.01095
- [26] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. 2023. Latent Consistency Models: Synthesizing High-Resolution Images with Few-Step Inference. arXiv:2310.04378
- [27] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. 2023. Latent-NeRF for Shape-Guided Generation of 3D Shapes and Textures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 12663–12673.
- [28] 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.
- [29] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2022. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. arXiv:2112.10741
- [30] 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:2212.08751
- [31] Alexander Quinn Nichol and Prafulla Dhariwal. 2021. Improved Denoising Diffusion Probabilistic Models. In *Proceedings of the 38th International Conference on Machine Learning*, Vol. 139. 8162–8171.
- [32] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. 2022. DreamFusion: Text-to-3D using 2D Diffusion. arXiv (2022).
- [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. arXiv:2103.00020
- [34] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical Text-Conditional Image Generation with CLIP Latents. arXiv:2204.06125
- [35] 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 (CVPR)*. 10684–10695.
- [36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation. arXiv:1505.04597
- [37] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. 2022. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. In *Advances in Neural Information Processing Systems*, Vol. 35. 36479–36494.
- [38] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Saxena, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. In *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Eds.), Vol. 35. 25278–25294.
- [39] Kihyuk Sohn, Honglak Lee, and Kinchen Yan. 2015. Learning Structured Output Representation using Deep Conditional Generative Models. In *Advances in Neural Information Processing Systems*, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (Eds.), Vol. 28.
- [40] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. Denoising Diffusion Implicit Models. arXiv:2010.02502 (2020).
- [41] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. 2023. Consistency Models. arXiv:2303.01469
- [42] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. 2020. Score-Based Generative Modeling through Stochastic Differential Equations. arXiv:2011.13456
- [43] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. 2024. DreamGaussian: Generative Gaussian Splatting for Efficient 3D Content Creation. arXiv:2309.16653
- [44] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Greg Shakhnarovich. 2023. Score Jacobian Chaining: Lifting Pretrained 2D Diffusion Models for 3D Generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 12619–12629.
- [45] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. 2023. ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation. In *Advances in Neural Information Processing Systems*, A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (Eds.), Vol. 36. 8406–8441.
- [46] Taoran Yi, Jiemin Fang, Junjie Wang, Guanjun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and Xinggang Wang. 2023. GaussianDreamer: Fast Generation from Text to 3D Gaussians by Bridging 2D and 3D Diffusion Models. arXiv:2310.08529
- [47] Xin Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Song-Hai Zhang, and Xiaojuan Qi. 2023. Text-to-3d with classifier score distillation. arXiv preprint arXiv:2310.19415 (2023).
- [48] Qingsheng Zhang, Molei Tao, and Yongxin Chen. 2023. gDDIM: Generalized denoising diffusion implicit models. arXiv:2206.05564
- [49] Junzhe Zhu, Peiye Zhuang, and Sanmi Koyejo. 2024. HiFA: High-fidelity Text-to-3D Generation with Advanced Diffusion Guidance. arXiv:2305.18766

929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044

1045 Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

1046	1103
1047	1104
1048	1105
1049	1106
1050	1107
1051	1108
1052	1109
1053	1110
1054	1111
1055	1112
1056	1113
1057	1114
1058	1115
1059	1116
1060	1117
1061	1118
1062	1119
1063	1120
1064	1121
1065	1122
1066	1123
1067	1124
1068	1125
1069	1126
1070	1127
1071	1128
1072	1129
1073	1130
1074	1131
1075	1132
1076	1133
1077	1134
1078	1135
1079	1136
1080	1137
1081	1138
1082	1139
1083	1140
1084	1141
1085	1142
1086	1143
1087	1144
1088	1145
1089	1146
1090	1147
1091	1148
1092	1149
1093	1150
1094	1151
1095	1152
1096	1153
1097	1154
1098	1155
1099	1156
1100	1157
1101	1158
1102	1159
	1160