# AnchorDS: Anchoring Dynamic Sources for Semantically Consistent Text-to-3D Generation

**Jiayin Zhu[1], Linlin Yang[2], Yicong Li[1*], Angela Yao[1],**

[1]National University of Singapore
[2]Communication University of China
zhujiayin@u.nus.edu, lyang@cuc.edu.cn, liyicong@u.nus.edu, ayao@comp.nus.edu.sg

## Abstract

Optimization-based text-to-3D methods distill guidance from 2D generative models via Score Distillation Sampling (SDS), but implicitly treat this guidance as static. This work shows that ignoring source dynamics yields inconsistent trajectories that suppress or merge semantic cues, leading to "semantic over-smoothing" artifacts. As such, we reformulate text-to-3D optimization as mapping a *dynamically evolving source* distribution to a fixed target distribution. We cast the problem into a dual-conditioned latent space, conditioned on both the text prompt and the intermediately rendered image. Given this joint setup, we observe that the image condition naturally anchors the current source distribution. Building on this insight, we introduce AnchorDS, an improved score distillation mechanism that provides state-anchored guidance with image conditions and stabilizes generation. We further penalize erroneous source estimates and design a lightweight filter strategy and fine-tuning strategy that refines the anchor with negligible overhead. AnchorDS produces finer-grained detail, more natural colours, and stronger semantic consistency, particularly for complex prompts, while maintaining efficiency. Extensive experiments show that our method surpasses previous methods in both quality and efficiency.

**Code** — https://github.com/viridityzhu/AnchorDS

## 1 Introduction

With the growing demand for 3D content creation in gaming and virtual reality, text-to-3D generation has emerged as a significant research frontier. One prominent solution that builds on this trend is Score Distillation Sampling (SDS) (Poole et al. 2023). SDS can leverage powerful off-the-shelf 2D diffusion models to guide optimization-based text-to-3D generation without large-scale 3D datasets.

SDS-guided text-to-3D generation can be interpreted as an optimization process that gradually shifts the distribution of rendered images from the current 3D representation (*i.e.*, the source distribution) toward a distribution defined by a text-conditioned diffusion model (*i.e.*, the target distribution) (McAllister et al. 2024). Under this interpretation, follow-up works have proposed enhancements follow two strategies. The first enhances the target distribution estimation by incorporating additional conditions, such as multiview images (Shi et al. 2023; Li et al. 2025) or depth
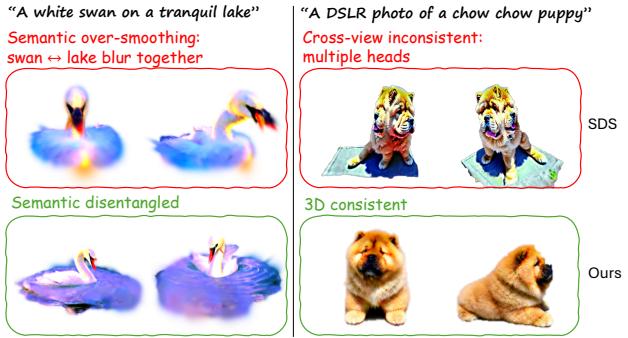
---

*Corresponding author.



Figure 1: **Comparison of Vanilla SDS vs. Ours.** SDS suffers from semantic over-smoothing (mixing swan and lake) and cross-view inconsistency (multiple heads). Ours achieves semantic disentanglement and 3D consistency across views.

maps (Qiu et al. 2024). The second addresses source distribution ismatch by refining the guidance, using reference-based scoring (Hertz, Aberman, and Cohen-Or 2023), enhanced negative prompts (McAllister et al. 2024), or variational optimization (Wang et al. 2023). Despite these advancements, SDS-based 3D generation still struggles with consistently preserving structural semantics across iterative updates. This manifests as two critical artifacts: (1) semantic over-smoothing, where object-specific features degenerate into homogenized, semantically ambiguous representations; and (2) cross-view inconsistency, where geometry and appearance are incoherent across perspectives (Fig. 1).

At its core, SDS optimization is formulated as a distribution transformation: it progressively shifts the rendered appearance of the 3D asset—from a *source distribution* reflecting the current 3D state toward a *target distribution* defined by a text-conditioned diffusion model. However, while the target distribution remains fixed to the text prompt, the source distribution is approximated using a static unconditional prior throughout optimization. This critically overlooks the *inherently non-stationary nature of the source distribution*: as the 3D representation evolves through optimization steps, rendered images continuously update, dynamically altering the source distribution. Consequently, SDS effectively discards accumulated structural information at each step by restarting from a semantics-agnostic prior. This fundamental mismatch between the static guidance mechanism and the evolving asset state destabilizes op-

timization, ultimately manifesting as the observed artifacts.

To address this issue, we propose to reframe score distillation as a dynamic editing process. Rather than treating generation as a one-shot projection from a static source, we view it as a progressive editing loop in which each optimization step refines the current 3D state based on the accumulated guidance from previous steps. Crucially, by explicitly recognizing and exploiting the evolving 3D state itself, our reformulation yields two advantages: 1) it supplies a stream of rich cues—geometry, colour, and semantics—that stabilise guidance and enforce cross-view consistency; and 2) feeding the state back into a pretrained conditional diffusion model enables dynamic, accurate, and lightweight source estimates without extra networks or handcrafted prompts.

Building on this perspective, we introduce AnchorDS, a dynamic form of SDS that anchors the source distribution using the rendered image at each optimization step. Specifically, we leverage a dual-conditioned diffusion model that incorporates both the text prompt and the intermediate image as guidance signals. Crucially, we observe that the diffusion model's predicted noise conditioned on the current rendering naturally encodes structural and semantic cues from the image condition. This noise prediction inherently anchors the source distribution by correlating the gradient guidance with the evolving 3D state, mitigating distribution drift without explicit constraints. Notably, the image condition does not constrain the target output directly, but serves as a contextual anchor that steers the generation. This design is not sensitive to the selection of the dual-conditioned model. Our method is robust when utilizing various popular image-conditioned adapters, *e.g*, IP-Adapter (Ye et al. 2023) and ControlNet (Zhang, Rao, and Agrawala 2023).

Despite this flexibility, anchoring dynamic sources is still non-trivial because the implicit latent space of the diffusion model exhibits a distribution mismatch with rendered images. To bridge the gap, we incorporate two complementary components. First, we reconstruct a pseudo-source image at each step, offering a metric for evaluating the quality of the estimated source distribution. Second, we introduce two practical enhancements: a simple yet effective *Filtering* mechanism that discards unreliable source predictions, and a lightweight *Fine-tuning* strategy that better aligns the diffusion model with the domain of rendered images.

Our contributions are summarized as follows:

1. We reveal the evolving nature of the source distribution in SDS and identify it as the root cause of semantic oversmoothing and inconsistent optimization trajectories.

2. We propose AnchorDS, a novel score distillation framework that dynamically anchors the source estimation via a dual-conditioned diffusion model. We further introduce a filtering mechanism and a lightweight fine-tuning strategy to better regularize the evolving source distribution.

3. Extensive experiments on T$^3$Bench (He et al. 2023) and a representative suite of challenging prompts demonstrate that our method outperforms state-of-the-art (SoTA) SDS variants in both generation quality and efficiency.

## 2   Related Works

**Text-Guided 3D Generation.** DreamFusion demonstrates that a NeRF can be trained from scratch by following Score Distillation Sampling (SDS) gradients derived from 2D diffusions (Poole et al. 2023). Subsequent works keep the same optimization-from-noise paradigm while improving efficiency or fidelity (Lin et al. 2023; Tang et al. 2024b; Yi et al. 2024a; Li et al. 2024). Although these methods deliver the highest visual quality, they remain slow and are still vulnerable to view inconsistency and Janus problems. To eliminate per-scene optimization, another line of work trains generators that map text directly to 3D latents (Jun and Nichol 2023; Nichol et al. 2022; Siddiqui et al. 2024; Xiang et al. 2024). While generation is fast, training quality rely on large-scale 3D datasets (Deitke et al. 2022; Fu et al. 2021; Collins et al. 2022), and the outputs still trail optimization methods in geometric and color accuracy (He et al. 2023). Hybrid strategies therefore combine a feed-forward initialization with subsequent Gaussian/NeRF refinement (Liang et al. 2024; Yi et al. 2024b,a).

**Conditional and Controllable Text-to-3D.** Extending the diffusion prior with extra modalities improves controllability. ControlNet (Zhang, Rao, and Agrawala 2023) has been adopted for depth, normal, or multi-view constraints (Huang et al. 2024; Li et al. 2025). IP-Adapter (Ye et al. 2023) let users steer style or identity with a reference image and have been plugged into 3D pipelines (Zeng et al. 2023). All these works, however, treat the additional image as an external positive condition supplied a priori. We instead employs the intermediate rendering itself as a dynamic *source* anchor supplying self-consistent guidance.

**Refining Score Distillation Sampling.** Recent work revisits SDS from theoretical and practical perspectives. (Yu et al. 2023; Tang et al. 2024a; Katzir et al. 2024) separate mode-seeking and variance terms to stabilize optimization. (Liang et al. 2024; Lukoianov et al. 2024) avoid first-order errors through DDIM sampling and inversion. Mismatch remedies attempt to align the diffusion prior with the 3D asset. DDS pairs each original image with a reference prompt (Hertz, Aberman, and Cohen-Or 2023), SDS-Bridge introduces a handcrafted prompt describing the poor 3D state (McAllister et al. 2024), and ProlificDreamer trains a LoRA branch to approximate the particle distribution (Wang et al. 2023). These solutions improve stability yet rely on static prompts, specific references, or auxiliary networks, which introduce new bias and overhead. Instead, we remove this dependency by directly feeding the current rendering into the diffusion prior, yielding faithful and bias-free guidance efficiently.

## 3   Analysis on the Issue of Source Distribution Estimation

We begin by revisiting the formulation of text-to-3D generation via Score Distillation Sampling (SDS) (Poole et al. 2023), and highlight its key limitation - a lack of awareness of the rendered appearance from the current 3D state. We then draw connections between SDS and 2D editing paradigms, and reinterpret SDS as a dynamic editing process that conditions on the evolving source.
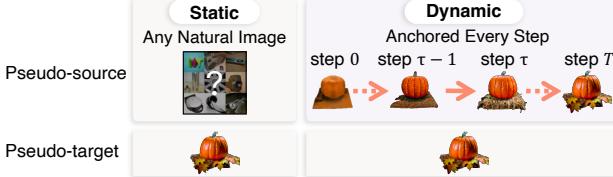
Figure 2: **Static vs. Dynamic Pseudo-Editing.** Static source estimation leverages an unconditional prior, misaligned with the actual 3D state. Instead, the dynamic pseudo-source reflects the evolving 3D renderings, ensuring faithful guidance.

## 3.1 Preliminaries

**Score Distillation Sampling (SDS).** SDS leverages pretrained 2D diffusion models to optimize the 3D generation. Specifically, it applies the denoising process of a 2D diffusion model to a rendered image from the 3D model, through which it distills a prior on the generated 3D output. Given a sampled noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ and a latent representation $z$ of the rendered image, the noisy latent at timestep $t$ is given by

$$z_t = \sqrt{\bar{\alpha}_t}\, z + \sqrt{1 - \bar{\alpha}_t}\, \epsilon. \quad (1)$$

The corresponding noisy image is then used to compute the SDS gradient $\nabla_\Theta \mathcal{L}_{SDS}(\phi, z)$, where the predicted noise is assumed to approximate the score function. Concretely, following the score matching view, the diffusion model $\phi$'s noise prediction $\hat{\epsilon}_\phi$ is defined as:

$$\hat{\epsilon}_\phi = -\sigma_t \nabla_{z_t} \log p(z_t; t, c), \quad \text{with } \sigma_t = \sqrt{1 - \bar{\alpha}_t}, \quad (2)$$

where $c$ denotes conditioning information (typically a text prompt $y$, or a combination $c = \{y, I\}$ when an image condition $I$ is included). This formulation directs the optimization toward regions of high density in the conditional distribution $p(z_t; t, c)$. Applied to a 3D model parameterized by $\Theta$, the gradient is given as:

$$\nabla_\Theta \mathcal{L}_{SDS}(\phi, z) = w(t) \left( \hat{\epsilon}_\phi^{CFG}(z_t; t, c) - \epsilon \right) \frac{\partial z_t}{\partial \Theta}, \quad (3)$$

where $\hat{\epsilon}_\phi^{CFG}(z_t; t, c)$ is the noise estimate (Eq. 2) under Classifier-Free Guidance (CFG).

**Classifier-Free Guidance (CFG).** CFG (Ho and Salimans 2022) balances conditional and unconditional predictions through adjustable weights. For a single-condition (text-only) setup, the CFG prediction is:

$$\hat{\epsilon}_\phi^{CFG}(z_t, t, y) = (1 + \omega)\, \hat{\epsilon}_\phi(z_t, t, y) - \omega\, \hat{\epsilon}_\phi(z_t, t, \emptyset), \quad (4)$$

where $\omega$ controls the strength of the guidance.

## 3.2 The Issue of Static Source Estimation

We first analyze why static source modeling in SDS causes artifacts (Fig. 1). Substituting Eq. 4 into Eq. 3 leads to two key terms $m_1$ and $m_2$:

$$\hat{\epsilon}_\phi^{CFG}(z_t; t, c) - \epsilon = (\omega - 1) \cdot \underbrace{\left( \hat{\epsilon}_\phi(z_t, t, c) - \hat{\epsilon}_\phi(z_t, t, \emptyset) \right)}_{m_1}$$
$$+ \underbrace{\hat{\epsilon}_\phi(z_t, t, \emptyset) - \epsilon}_{m_2}. \quad (5)$$

A sufficiently large scaling of $\omega$, *e.g.*, $\omega = 100$ will cause $m_1$ to dominate the guidance, while $m_2$ reduces variance of $\hat{\epsilon}_\phi^{CFG}$ (McAllister et al. 2024; Tang et al. 2024a). Interpret-

ing $m_1$ as a score function:

$$m_1 = -\sigma_t \left( \nabla_{z_t} \log p(z_t; t, c) - \nabla_{z_t} \log p(z_t; t, \emptyset) \right), \quad (6)$$

and we see SDS pushes samples toward the higher-density regions of the conditional distribution $p(z_t; t, c)$ while moving away from the unconditional prior $p(z_t; t, \emptyset)$.

This mirrors 2D image editing (*e.g*, DDIB (Su et al. 2023), SDEdit (Meng et al. 2022)), where edits aim to find the optimal mapping from a source image to a target image as two distributions. However, SDS approximates the source distribution—which should reflect the current 3D state—using the unconditional prior $p(z_t; t, \emptyset)$. The limitation becomes clear when rearranging Eq. 6 into a pseudo-editing formulation. Specifically, we invert Eq. 1 to define the pseudo-target and pseudo-source latent reconstructions:

$$\hat{z}_{t\rightarrow 0}^{\text{target}} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( z_t - \sqrt{1 - \bar{\alpha}_t}\, \hat{\epsilon}_\phi(z_t, t, y) \right),$$
$$\hat{z}_{t\rightarrow 0}^{\text{source}} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( z_t - \sqrt{1 - \bar{\alpha}_t}\, \hat{\epsilon}_\phi(z_t, t, \emptyset) \right). \quad (7)$$

Then, the core update term simplifies elegantly as:

$$m_1 = \eta \left( \hat{z}_{t\rightarrow 0}^{\text{target}} - \hat{z}_{t\rightarrow 0}^{\text{source}} \right), \quad \text{where} \quad \eta = \frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}}. \quad (8)$$

Examining $\hat{z}_{t\rightarrow 0}^{\text{source}}$, we reveal that it loses critical information about the original rendering: (1) $z_t$ is a noise-corrupted version that discards semantic content, while (2) $\hat{\epsilon}_\phi(z_t, t, \emptyset) = -\sigma_t \nabla_{z_t} \log p(z_t; t, \emptyset)$ pushes toward an unconditional image prior that averages over diverse natural images. Crucially, neither term encodes the evolving 3D asset's current state, forcing $\hat{z}_{t\rightarrow 0}^{\text{source}}$ to inherit this static, averaged characteristic rather than reflecting the dynamic source distribution. As a result, it fails to capture the semantics of intermediate renderings. The resulting update directions become inconsistent with the 3D asset's actual appearance, causing oversmoothing and cross-view inconsistency in Fig. 1.

Rather than static, we argue that the source should evolve with the 3D model (Fig. 2). We therefore reinterpret the text-to-3D optimization as a *dynamic edit*: at each step, we refine the 3D asset by (1) preserving favorable attributes (*e.g*, structural semantics consistent with current state) via faithful source modeling, and (2) correcting undesirable features (*e.g*, deviations from target distribution) through prompt alignment. This resolves the inconsistency in Fig. 1 by anchoring updates to the actual state rather than a static prior.

# 4 Method

## 4.1 Dynamic Evolution of the Source Distribution

We formalize a critical insight regarding SDS-based text-to-3D generation: the source image distribution $P_\theta(x)$ evolves dynamically throughout optimization rather than remaining static, as visualized in Fig. 3b. Mathematically, we reframe SDS optimization as a progressive transformation process:

$$P_\theta^{(0)}(x) \xrightarrow{\nabla_\Theta \mathcal{L}_{SDS}} P_\theta^{(1)}(x) \rightarrow \cdots \rightarrow P_\theta^{(T)}(x), \quad (9)$$

where $P_\theta^{(\tau)}(x)$ denotes the time-varying source distribution at optimization step $\tau$. Initially, $P_\theta^{(0)}(x)$ resembles a diffuse unconditional prior when the 3D model $\Theta$ is randomly initialized. At each iteration, the gradient update $\nabla_\Theta \mathcal{L}_{SDS}$
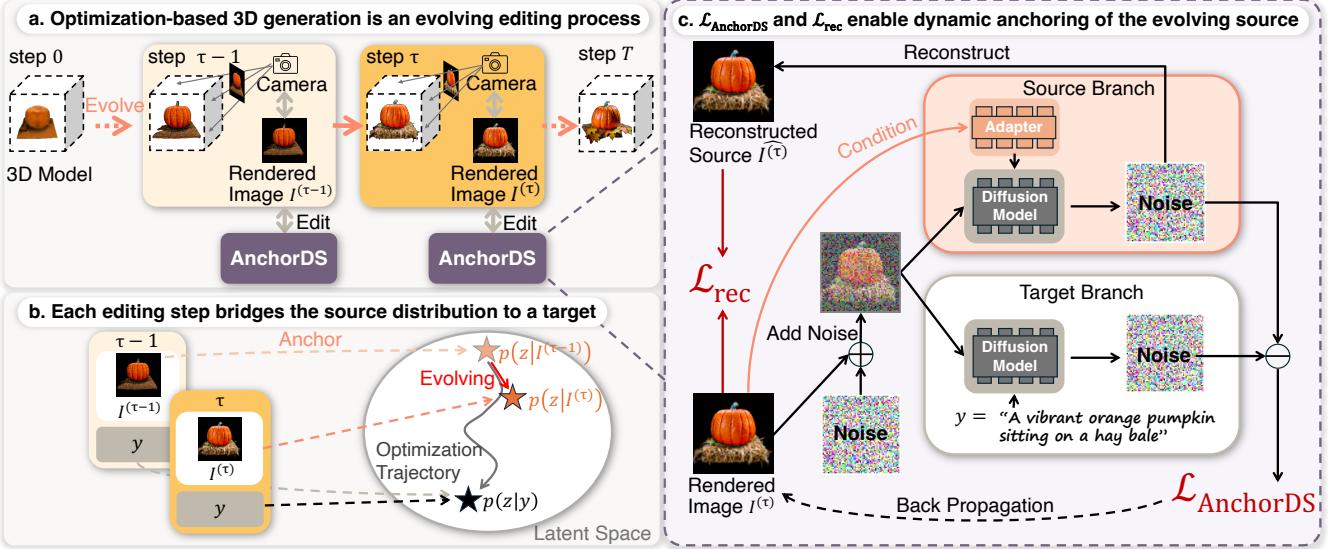
Figure 3: **Overview of Our AnchorDS for Text-to-3D Generation. (a)** Optimization-based 3D generation as an evolving editing process where each step refines the 3D model guided by AnchorDS. **(b)** Each editing step bridges the source distribution to a target distribution through dynamic anchoring in latent space. **(c)** Technical framework showing how $L_{AnchorDS}$ and $L_{rec}$ enable dynamic source anchoring, with source and target branches processing the evolving 3D content through pretrained diffusion models.

transports probability mass such that:

$$P_\theta^{(\tau+1)}(x) = \mathcal{T}\left(P_\theta^{(\tau)}(x), \nabla_\Theta \mathcal{L}_{\mathrm{SDS}}(z, y, \hat{\boldsymbol{\epsilon}}_t)\right), \quad (10)$$

where $\mathcal{T}(\cdot)$ represents the distributional shift operator. This process progressively morphs $P_\theta^{(\tau)}(x)$ toward the target distribution $P_{\mathrm{target}}(x \mid y)$ defined by the text prompt $y$. The continuous distributional evolution stems from the mass transport interpretation of score-based dynamics (De Bortoli et al. 2021), establishing this formulation's generality across score-distillation-based 3D generation methods.

Next, we introduce our AnchorDS to ensure accurate source distribution estimation at every optimization step.

## 4.2 Score Distillation via Dynamic Source Anchoring

AnchorDS dynamically anchors the score guidance at the evolving 3D state through image conditioning. Formally, at optimization step $\tau$ with rendered view $I^{(\tau)}$, we compute the guidance gradient as:

$$g_t^{(\tau)} = \hat{\boldsymbol{\epsilon}}_\phi(z_t; t, y) - \hat{\boldsymbol{\epsilon}}_\phi(z_t; t, \emptyset, I^{(\tau)}), \quad (11)$$

where $\hat{\boldsymbol{\epsilon}}_\phi(z_t; t, y)$ targets the text-conditioned distribution, while $\hat{\boldsymbol{\epsilon}}_\phi(z_t; t, \emptyset, I^{(\tau)})$ anchors the current source distribution $P_\theta^{(\tau)}(x)$. This differential formulation directs updates from the current state toward the target distribution.

Essentially, incorporating an image condition $I^{(\tau)}$ recasts the problem into a dual-conditioned latent space. This preserves the text-conditioned target's effectiveness while leveraging image conditions to anchor the source distribution. Specifically, $I^{(\tau)}$ anchors the current source by providing structural grounding without content constraints, enabling contextual editing rather than output restriction. It yields two principal advantages: (1) Conditioning with $I^{(\tau)}$ at continuous $\tau$ maintains alignment between guidance and 3D state, mitigating drift and oversmoothing. (2) Source an-

choring requires only a single additional U-Net forward pass per iteration, processed in parallel with the original pass, thus maintaining identical runtime to standard SDS.

AnchorDS bridges 2D editing principles with 3D generation through a key insight: pretrained diffusion models inherently possess image inversion capabilities within their conditional architecture. While 2D editing requires explicit inversion techniques, AnchorDS elegantly leverages this intrinsic property—directly utilizing the model's natural ability to map images to latent distributions, achieving precise source anchoring without auxiliary inversion costs.

## 4.3 Source Anchoring via Image Conditioning

The effectiveness of AnchorDS depends critically on an accurate estimation of the current source distribution via image-conditioned diffusion models. We denote by $I^{(\tau)} = R(\mathcal{X}_{3D}^{(\tau)})$ the rendered image of the 3D representation $\mathcal{X}_{3D}^{(\tau)}$.
**Choice of Image Condition.** In general, one could apply a preprocessing function $E(\cdot)$ to obtain a conditioning signal $\bar{I}^{(\tau)} = E(I^{(\tau)})$ (e.g, converting $I^{(\tau)}$ to a normal map, depth map, etc.) The conditioning signal must retain essential structural and semantic information from the 3D rendering while eliminating irrelevant noise. We initially consider that the identity mapping $E(I^{(\tau)}) = I^{(\tau)}$ is particularly effective as it preserves maximal information about the current state. This aligns with evidence (Kadosh et al. 2025) that image-conditioned diffusion models learn invertible mappings between images and noise—crucial for source estimation. Meanwhile, we find that our method remains compatible with alternative signals (e.g, normal maps) provided that they sufficiently capture the 3D content's core attributes.
**Pseudo-Source Reconstruction.** With the image-conditioned model in place, we can obtain an explicit estimate of the current source image distribution at any diffusion step. Given a noisy latent $z_t$ (at noise level $t$) that
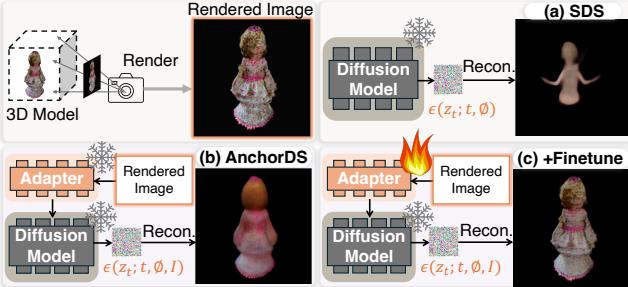
Figure 4: **Source Estimation Strategies.** (a) Vanilla SDS poorly reconstructs the source. (b) AnchorDS, conditioned on the current rendering, recovers geometry and color faithfully. (c) Our fine-tuning strategy achieves an accurate match to the source image.

was produced from the current image $I^{(\tau)}$, the model's image conditioned prediction $\hat{\epsilon}_\phi(z_t; t, \emptyset, I^{(\tau)})$ allows us to reconstruct a pseudo-reconstructed source image:

$$\hat{z}_{t\to0}^{\text{anchored},(\tau)} = \frac{1}{\sqrt{\bar{\alpha}_t}}\left(z_t - \sqrt{1 - \bar{\alpha}_t}\,\hat{\epsilon}_\phi(z_t; t, \emptyset, I^{(\tau)})\right), \quad (12)$$

which is the model's one-step estimate of the denoised latent at timestep $t$ – essentially a guess of the original image $I^{(\tau)}$ (latent) before noise was added. With an image decoder $\varepsilon(\cdot)$, we can now explicitly evaluate reconstruction accuracy as:

$$\mathcal{L}_{\text{rec}} = \left\|\varepsilon(\hat{z}_{t\to0}^{\text{anchored},(\tau)}) - I^{(\tau)}\right\|_2^2. \quad (13)$$

This reconstruction not only provides a direct metric for source estimation quality but also enables two complementary mechanisms: (1) filtering out unstable predictions, and (2) fine-tuning the image adapter to better align with rendered images—both crucial for stable and accurate 3D optimization. We implement two strategies discussed below.

**Filtering.** A straightforward way to employ the objective (Eq. 13) is to apply a threshold-based filter to exclude unreliable source estimations:

$$\mathcal{M}_{\text{Filter}} = \begin{cases} 1 & \text{if } \mathcal{L}_{\text{rec}} < \gamma \\ 0 & \text{otherwise} \end{cases} \quad (14)$$

$$\mathcal{L}_{\text{AnchorDS}} = \mathcal{M}_{\text{Filter}} \cdot \mathcal{L}_{\text{AnchorDS}}.$$

Here, $\mathcal{L}_{\text{rec}}$ acts as a stabilizer, filtering out any spurious deviations in the anchored source prediction. Empirically, this translates to improved stability (no sudden jarring updates) and better preservation of existing content.

**Finetuning Image Adapter.** Another approach is to apply $\mathcal{L}_{\text{rec}}$ to fine-tune the image adapter. Although pretrained 2D models are powerful enough to model the real-world image distribution, there remains a gap between the real-world and the synthesized rendered image distribution. Intuitively, we aim to let 2D models "see" the real data. We address this through lightweight fine-tuning of the image adapter using $\mathcal{L}_{\text{rec}}$. This fine-tuning is minimal—only unfreezing one single layer of the image adapter is sufficient (increasing optimization time from ~25 minutes to ~30 minutes using 3DGS pipeline (Yi et al. 2024a)), while significantly improving source estimation accuracy (see Fig. 4c).

Minimizing $\mathcal{L}_{\text{rec}}$ forces $\hat{\epsilon}_\phi(z_t; t, \emptyset, I^{(\tau)})$ to become consistent with the current image $I^{(\tau)}$. Note $\hat{\epsilon}_\phi(z_t; t, \emptyset, I^{(\tau)})$ is *not* trained to match sampled noise; Instead, its role is to "invert" the current image into latent space. Eq. 12 sim-

ply uses this predicted noise to retrieve the model's internal guess of the clean image $I^{(\tau)}$, and $\mathcal{L}_{\text{rec}}$ ties that guess further. With this enhanced source estimate term, guidance $g_t^{(\tau)}$ (Eq. 11) accurately reflects the difference between distributions of images that ensembles $I^{(\tau)}$ and those that fulfill $y$.

## 4.4 Implementation

**Choice of Diffusion Model.** To incorporate $\bar{I}^{(\tau)}$ as a condition, we leverage existing pre-trained dual-conditional diffusion models. A straightforward option is an image-to-image diffusion model such as IP2P (Brooks, Holynski, and Efros 2023). However, IP2P is fine-tuned for image editing and we observed it gives inconsistent guidance under the high CFG weights required for SDS; in practice, using IP2P led to unstable and desaturated results for text-to-3D, especially at large CFG scales. ControlNet (Zhang, Rao, and Agrawala 2023) processes derived maps (e.g, normal maps, sketches) but lacks direct training on raw images. Conversely, IP-Adapter (Ye et al. 2023) conditions Stable Diffusion (1.5) on unaltered images via an auxiliary latent, preserving the model's expressive power without constraining image content. Crucially, our framework generalizes across adapters. We adopt IP-Adapter for primary experiments due to its direct image conditioning, while adopting ControlNet shows comparative results (detailed in Sec. 5).

**Pipeline.** Given a text prompt $y$ and an optional initial 3D representation (which could be random or from a rough generator), our pipeline operates iteratively as illustrated in Fig. 3c. At each optimization step $\tau$, we render the current 3D model from a random viewpoint to obtain image $I^{(\tau)}$, then encode it into the diffusion latent space and add noise at a random timestep $t$ to produce $z_t$. We then query the diffusion model twice: once with text condition $y$ to obtain the target prediction $\hat{\epsilon}_\phi(z_t; t, y)$, and once with an empty text and $I^{(\tau)}$ to obtain the source prediction $\hat{\epsilon}_\phi(z_t; t, \emptyset, I^{(\tau)})$. The obtained AnchorDS guidance (Eq. 11) is then backpropagated through the rendering pipeline to update the 3D parameters. For the fine-tuning variant, we periodically update the unfrozen layer of the image adapter using $\mathcal{L}_{\text{rec}}$.

## 5 Experiments

**Experimental Setup.** We aim to evaluate the effectiveness of AnchorDS in addressing the limitations of SDS guidance, specifically source distribution mismatch and semantic over-smoothing. We compare against vanilla SDS (Poole et al. 2023) and two representative methods tackling similar issues: SDS-Bridge (McAllister et al. 2024) and ProlificDreamer (VSD)(Wang et al. 2023). To validate robustness, we test across both 3D Gaussian Splatting (3DGS)(Yi et al. 2024a) and NeRF-based pipelines.

**Evaluation Metrics.** Following (McAllister et al. 2024; Lee, Sohn, and Shin 2024; Dong et al. 2024), we employ CLIP similarity (Radford et al. 2021) between text prompts and rendered images to assess generation alignment. Additionally, we adopt T$^3$Bench (He et al. 2023)'s quality metric, evaluating the 3D visual quality using pre-trained language and visual models. Since the target estimation term in Eq. 6 remains unchanged, text alignment is maintained by design, letting us focus primarily on quality improvements.
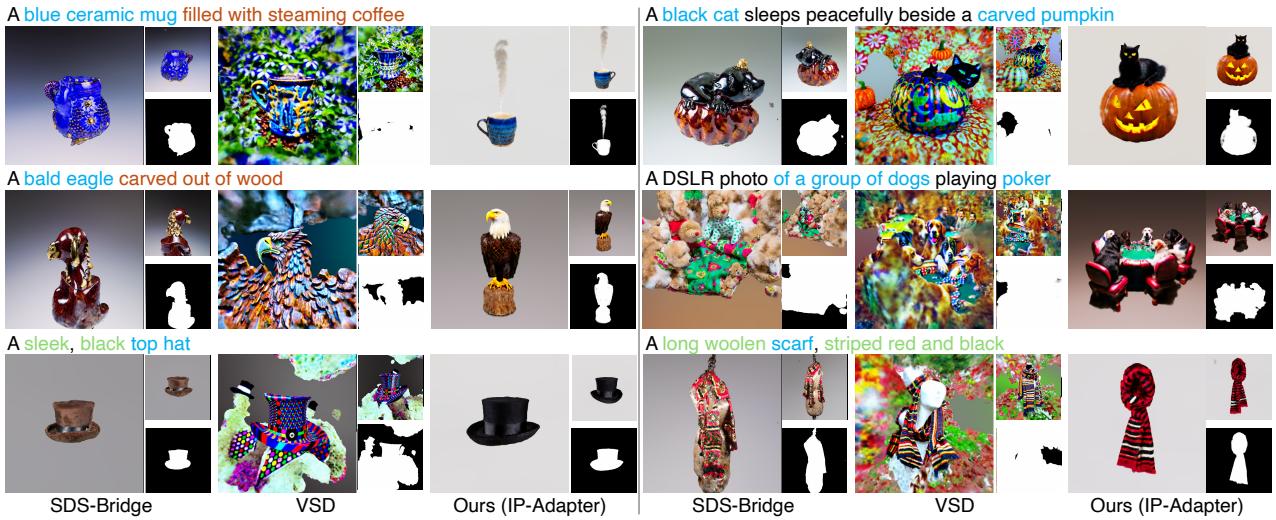
A blue ceramic mug filled with steaming coffee

A black cat sleeps peacefully beside a carved pumpkin

A bald eagle carved out of wood

A DSLR photo of a group of dogs playing poker

A sleek, black top hat

A long woolen scarf, striped red and black

| SDS-Bridge | VSD | Ours (IP-Adapter) | SDS-Bridge | VSD | Ours (IP-Adapter) |

Figure 5: **Qualitative Comparison to Existing Methods on SD 1.5** on complex text prompts, including (1) objects with rich or mixed semantics; (2) *multiple* distinct objects; (3) objects with fine-grained details. SDS-Bridge produces biased material textures and inaccurate semantics, likely due to new biases introduced by the handcrafted prompt. VSD fails to recover coherent structures and exhibits exaggerated, unrealistic colours. In contrast, our method consistently generates semantically faithful and structurally accurate results.

Table 1: **Quantitative Comparison.** Q1–Q3 report averaged ranking of each method.

| Method | Base Model | CLIP ↑ | Q1 ↓ | Q2 ↓ | Q3 ↓ |
|---|---|---|---|---|---|
| VSD | SD 2.1 | 0.352 | 1.84 | 1.85 | 1.79 |
| Ours (ControlNet) | | **0.369** | **1.16** | **1.15** | **1.21** |
| VSD | SD 1.5 | 0.281 | 1.99 | 2.00 | 2.08 |
| SDS-Bridge | | 0.233 | 2.38 | 2.35 | 2.29 |
| Ours (IP-Adapter) | | **0.334** | **1.63** | **1.66** | **1.63** |

Q1: Which one has the best 3D consistency?
Q2: Which one shows accurately what the text describes?
Q3: Which one looks most realistic and natural?

## 5.1 Comparison with State-of-the-Art

To systematically evaluate SDS's limitations, we curate 50 complex prompts from previous works (He et al. 2023; Poole et al. 2023), covering three challenging categories: fine-grained details, rich/mixed semantics, and multiple-object compositions. This evaluation scale aligns with established practices (McAllister et al. 2024; Liang et al. 2024; Lee, Sohn, and Shin 2024; Dong et al. 2024). We adopt the NeRF-based pipeline (Wang et al. 2023) across all experiments. While VSD demonstrates sensitivity to base models and typically performs reasonably only on Stable Diffusion (SD) 2.1, we evaluate on both SD 1.5 and SD 2.1. Our method (AnchorDS with Finetuning) uses IP-Adapter on SD 1.5 and ControlNet on SD 2.1 as image conditioners.

**Human Evaluation.** We also conduct a comprehensive human preference evaluation on Amazon Mechanical Turk to assess methods across multiple dimensions. The evaluation comprises 20 batches covering our 50 complex prompts evaluated by 912 unique participants. For each comparison, evaluators rank all methods according to three questions: 3D consistency, text alignment, and visual quality, respectively.

**Results.** Table 1 shows our method consistently outperforms all methods in CLIP similarity and human preference evaluations. Qualitative comparisons in Fig. 5 for SD 1.5 as base model, and Fig. 6 for SD 2.1 as base model, and a multiview structural consistency comparison Fig. 7 also demonstrate our consistent generation of semantically faithful and struc-

Table 2: **Quantitative Comparison with Baseline on T$^3$Bench Quality Metric.** AnchorDS incorporates image conditioning via IP-Adapter or ControlNet. Filter and Finetune represent alternative strategies for enhanced source distribution estimation.

| Method | All↑ | Single↑ | Surr↑ | Multi↑ |
|---|---|---|---|---|
| SDS (DreamFusion) | 20.5 | 24.9 | 19.3 | 17.3 |
| SDS (GaussianDreamer) | 29.7 | 42.3 | 26.1 | 20.6 |
| AnchorDS (IP-Adapter) | 30.7 | 43.0 | 24.8 | 24.5 |
| + Filter | 32.8 | 44.1 | 27.9 | **26.5** |
| + Finetune | **33.3** | <u>45.3</u> | <u>29.0</u> | <u>25.7</u> |
| AnchorDS (ControlNet) | 30.8 | 43.9 | 27.2 | 21.3 |
| + Filter | <u>33.2</u> | **46.1** | **29.4** | 24.0 |
| + Finetune | 32.9 | 45.0 | 28.6 | 25.2 |

turally accurate results, while SDS-Bridge produces artifacts introduced by biased negative prompts, and VSD fails to accurately model the image distributions. Please refer to the Appendix for theoretical discussions.

## 5.2 Baseline Comparison

**SDS Baselines.** We comprehensively evaluate on T$^3$Bench (He et al. 2023), a comprehensive benchmark containing 300 prompts across single object generation, object with surrounding generation, and multiple object generation categories. We compare against both DreamFusion-SDS and GaussianDreamer-SDS baselines. To ensure fair evaluation, we maintain all pipeline components constant except for the score guidance, using the state-of-the-art Gaussian-Dreamer (Yi et al. 2024a) framework as our baseline. All methods are built upon SD 1.5.

**Comparison with SDS.** Quantitative results in Table 2 demonstrate that we consistently surpass SDS in all prompt categories. Our method achieves superior visual quality as illustrated in Fig. 8, with notable improvements in finer-grained detail, the preservation of complex semantic attributes, and the separation of multiple objects. The improvements align with our theoretical framework: Being aware of the evolving 3D state, AnchorDS generates along coherent semantic paths without mixing disparate attributes.

Figure 6: **Qualitative Comparison to Existing Methods on SD 2.1.** Our method surpasses VSD, consistently capturing complex semantics (even challenging elements like barns and brick walls) while ensuring superior 3D structural accuracy.
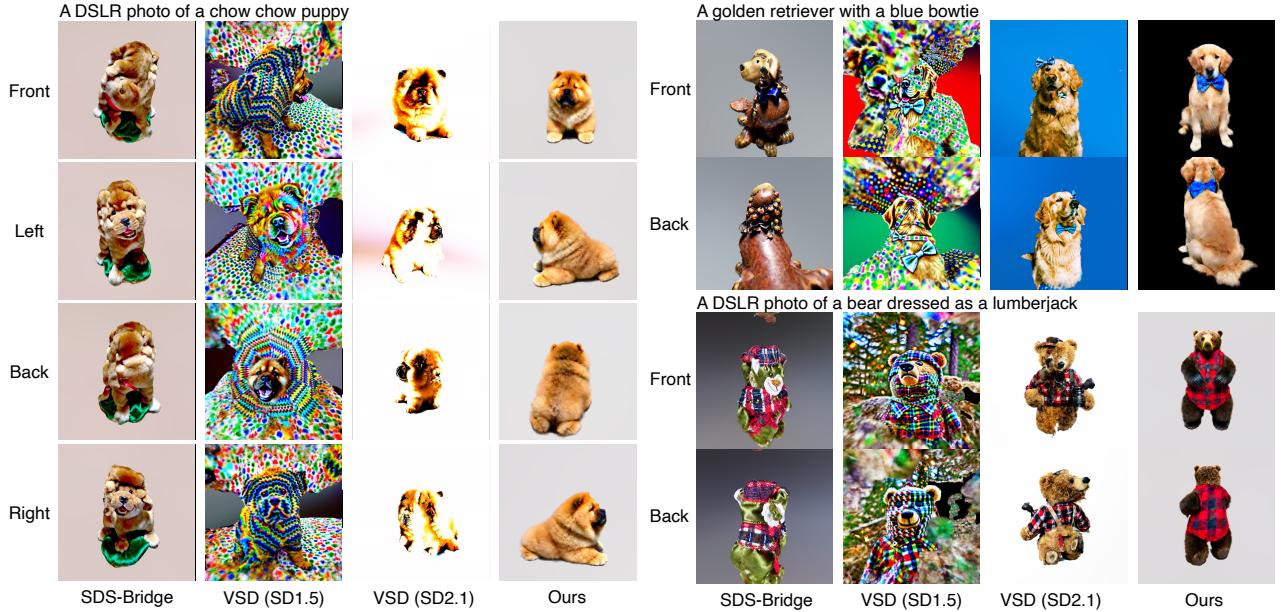


Figure 7: **Multi-view Comparison.** Ours demonstrates significantly better quality and multi-view consistency, despite not incorporating common 3D consistency strategies such as Perp-Neg (Armandpour et al. 2023). In contrast, although VSD introduces view direction as an additional condition, it still suffers from a severe Janus problem. SDS-Bridge fails to maintain consistent object semantics.
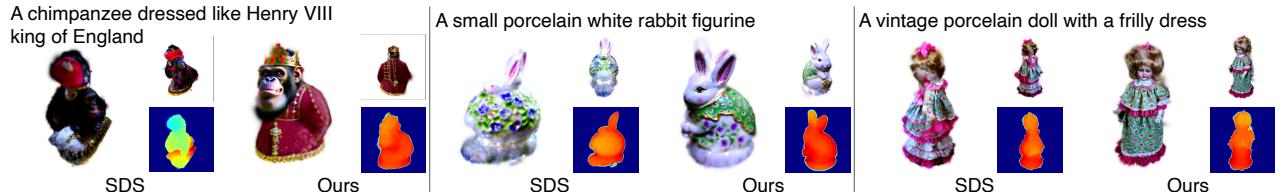


Figure 8: **Qualitative Comparison with Baseline (SDS).** Ours shows superior visual quality in finer-grained detail (*e.g*, doll's face). For prompts describing a single object with complex semantic details (*e.g*, porcelain white rabbit figurine), SDS tends to compress or mix parts of the information, whereas ours successfully preserves the full range of semantic attributes. For prompts involving multiple objects (*e.g*, chimpanzee and dresses), SDS appears mixed and blurry, while ours accurately separates distinct objects.

## 5.3 Ablation Studies

Quantitative ablation results for our key components are shown in Table 2. AnchorDS consistently outperforms vanilla SDS through source anchoring via image conditioning. Our Filtering and Finetuning strategies provide complementary approaches for enhancing source estimation accuracy, with Finetuning achieving optimal performance. Please refer to Appendix for more ablations across prior models.

## 6 Conclusion

We demonstrate that the source distribution dynamically evolves during text-to-3D optimization—a fundamental property that has been largely overlooked by existing meth-

ods. While prior efforts have focused on reducing trajectory estimation error or improving the guidance prior, they continue to exhibit critical issues such as semantic over-smoothing and multi-view inconsistency, due to inaccurate modeling of the evolving source. To address this, we introduce AnchorDS, which anchors the dynamic source distribution by casting the problem into a dual-conditioned latent space and conditioning on rendering images. Experimental results confirm our method effectively mitigates these issues and achieves superior fidelity and visual quality.

## References

Armandpour, M.; Sadeghian, A.; Zheng, H.; Sadeghian, A.; and Zhou, M. 2023. Re-imagine the negative prompt algo-

rithm: Transform 2d diffusion into 3d, alleviate janus problem and beyond. *arXiv preprint arXiv:2304.04968*. 7

Brooks, T.; Holynski, A.; and Efros, A. A. 2023. Instructpix2pix: Learning to follow image editing instructions. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 18392–18402. 5, 10

Collins, J.; Goel, S.; Deng, K.; Luthra, A.; Xu, L.; Gundogdu, E.; Zhang, X.; Vicente, T. F. Y.; Dideriksen, T.; Arora, H.; Guillaumin, M.; and Malik, J. 2022. ABO: Dataset and Benchmarks for Real-World 3D Object Understanding. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 21094–21104. IEEE. 2

De Bortoli, V.; Thornton, J.; Heng, J.; and Doucet, A. 2021. Diffusion schrödinger bridge with applications to score-based generative modeling. *Advances in neural information processing systems*, 34: 17695–17709. 4

Deitke, M.; Schwenk, D.; Salvador, J.; Weihs, L.; Michel, O.; VanderBilt, E.; Schmidt, L.; Ehsani, K.; Kembhavi, A.; and Farhadi, A. 2022. Objaverse: A Universe of Annotated 3D Objects. *arXiv preprint arXiv:2212.08051*. 2

Dong, W.; Yang, B.; Ma, L.; Liu, X.; Cui, L.; Bao, H.; Ma, Y.; and Cui, Z. 2024. Coin3d: Controllable and interactive 3d assets generation with proxy-guided conditioning. In *ACM SIGGRAPH Conf. Proc.*, 1–10. 5, 6

Fu, H.; Jia, R.; Gao, L.; Gong, M.; Zhao, B.; Maybank, S.; and Tao, D. 2021. 3D-FUTURE: 3D Furniture Shape with TextURE. *Int. J. Comput. Vis.*, 129(12): 3313–3337. 2

He, Y.; Bai, Y.; Lin, M.; Zhao, W.; Hu, Y.; Sheng, J.; Yi, R.; Li, J.; and Liu, Y.-J. 2023. T3Bench: Benchmarking Current Progress in Text-to-3D Generation. *arXiv preprint arXiv:2310.02977*. 2, 5, 6, 11, 12

Hertz, A.; Aberman, K.; and Cohen-Or, D. 2023. Delta Denoising Score. In *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2328–2337. 1, 2

Ho, J.; and Salimans, T. 2022. Classifier-Free Diffusion Guidance. *arXiv preprint arXiv:2207.12598*. 3

Huang, T.; Zeng, Y.; Zhang, Z.; Xu, W.; Xu, H.; Xu, S.; Lau, R. W.; and Zuo, W. 2024. Dreamcontrol: Control-based text-to-3d generation with 3d self-prior. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 5364–5373. 2

Jun, H.; and Nichol, A. 2023. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*. 2, 10, 12

Kadosh, E.; Goren, N.; Patashnik, O.; Garibi, D.; and Cohen-Or, D. 2025. Tight Inversion: Image-Conditioned Inversion for Real Image Editing. *arXiv preprint arXiv:2502.20376*. 4

Katzir, O.; Patashnik, O.; Cohen-Or, D.; and Lischinski, D. 2024. Noise-free Score Distillation. In *Proc. Int. Conf. Learn. Represent.* 2

Lee, K.; Sohn, K.; and Shin, J. 2024. DreamFlow: High-quality text-to-3D generation by Approximating Probability Flow. In *Proc. Int. Conf. Learn. Represent.* 5, 6

Li, W.; Chen, R.; Chen, X.; and Tan, P. 2024. Sweet-Dreamer: Aligning Geometric Priors in 2D diffusion for

Consistent Text-to-3D. In *Proc. Int. Conf. Learn. Represent.* 2

Li, Z.; Chen, Y.; Zhao, L.; and Liu, P. 2025. Controllable Text-to-3D Generation via Surface-Aligned Gaussian Splatting. In *Proc. Int. Conf. 3D Vis.* 1, 2

Liang, Y.; Yang, X.; Lin, J.; Li, H.; Xu, X.; and Chen, Y. 2024. Luciddreamer: Towards high-fidelity text-to-3d generation via interval score matching. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 6517–6526. 2, 6, 11

Lin, C.-H.; Gao, J.; Tang, L.; Takikawa, T.; Zeng, X.; Huang, X.; Kreis, K.; Fidler, S.; Liu, M.-Y.; and Lin, T.-Y. 2023. Magic3d: High-resolution text-to-3d content creation. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 300–309. 2

Lukoianov, A.; S'aez de Oc'ariz Borde, H.; Greenewald, K.; Guizilini, V.; Bagautdinov, T.; Sitzmann, V.; and Solomon, J. M. 2024. Score distillation via reparametrized ddim. *Adv. Neural Inf. Process. Syst.*, 37: 26011–26044. 2

McAllister, D.; Ge, S.; Huang, J.-B.; Jacobs, D. W.; Efros, A. A.; Holynski, A.; and Kanazawa, A. 2024. Rethinking Score Distillation as a Bridge Between Image Distributions. In *Proc. Int. Conf. Learn. Represent.* 1, 2, 3, 5, 6, 10

Meng, C.; He, Y.; Song, Y.; Song, J.; Wu, J.; Zhu, J.-Y.; and Ermon, S. 2022. SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations. In *Proc. Int. Conf. Learn. Represent.* 3

Nichol, A.; Jun, H.; Dhariwal, P.; Mishkin, P.; and Chen, M. 2022. Point-E: A System for Generating 3D Point Clouds from Complex Prompts. *ArXiv*, abs/2212.08751. 2

Poole, B.; Jain, A.; Barron, J. T.; and Mildenhall, B. 2023. DreamFusion: Text-to-3D using 2D Diffusion. In *Proc. Int. Conf. Learn. Represent.* 1, 2, 5, 6

Qiu, L.; Chen, G.; Gu, X.; Zuo, Q.; Xu, M.; Wu, Y.; Yuan, W.; Dong, Z.; Bo, L.; and Han, X. 2024. Richdreamer: A generalizable normal-depth diffusion model for detail richness in text-to-3d. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 9914–9925. 1

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *Proc. Int. Conf. Mach. Learn.*, 8748–8763. PMLR. 5

Shi, Y.; Wang, P.; Ye, J.; Long, M.; Li, K.; and Yang, X. 2023. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*. 1

Siddiqui, Y.; Monnier, T.; Kokkinos, F.; Kariya, M.; Kleiman, Y.; Garreau, E.; Gafni, O.; Neverova, N.; Vedaldi, A.; Shapovalov, R.; and Novotny, D. 2024. Meta 3D Asset-Gen: Text-to-Mesh Generation with High-Quality Geometry, Texture, and PBR Materials. In Globerson, A.; Mackey, L.; Belgrave, D.; Fan, A.; Paquet, U.; Tomczak, J.; and Zhang, C., eds., *Adv. Neural Inf. Process. Syst.*, volume 37, 9532–9564. Curran Associates, Inc. 2

Su, X.; Song, J.; Meng, C.; and Ermon, S. 2023. Dual Diffusion Implicit Bridges for Image-to-Image Translation. In *Proc. Int. Conf. Learn. Represent.* 3

Tang, B.; Wang, J.; Wu, Z.; and Zhang, L. 2024a. Stable Score Distillation for High-Quality 3D Generation. *arXiv preprint arXiv:2312.09305*, (arXiv:2312.09305). 2, 3, 11

Tang, J.; Ren, J.; Zhou, H.; Liu, Z.; and Zeng, G. 2024b. DreamGaussian: Generative Gaussian Splatting for Efficient 3D Content Creation. In *Proc. Int. Conf. Learn. Represent.* 2

Team, T. H. 2025. Hunyuan3D 2.1: From Images to High-Fidelity 3D Assets with Production-Ready PBR Material. arXiv:2506.15442. 12

Wang, Z.; Lu, C.; Wang, Y.; Bao, F.; LI, C.; Su, H.; and Zhu, J. 2023. ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation. In Oh, A.; Neumann, T.; Globerson, A.; Saenko, K.; Hardt, M.; and Levine, S., eds., *Adv. Neural Inf. Process. Syst.*, volume 36, 8406–8441. Curran Associates, Inc. 1, 2, 5, 6, 10, 11

Xiang, J.; Lv, Z.; Xu, S.; Deng, Y.; Wang, R.; Zhang, B.; Chen, D.; Tong, X.; and Yang, J. 2024. Structured 3D Latents for Scalable and Versatile 3D Generation. *arXiv preprint arXiv:2412.01506*. 2, 11, 12

Yan, R.; Wu, K.; and Ma, K. 2024. Flow Score Distillation for Diverse Text-to-3D Generation. *arXiv preprint arXiv:2405.10988*. 11

Ye, H.; Zhang, J.; Liu, S.; Han, X.; and Yang, W. 2023. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*. 2, 5, 10

Yi, T.; Fang, J.; Wang, J.; Wu, G.; Xie, L.; Zhang, X.; Liu, W.; Tian, Q.; and Wang, X. 2024a. Gaussiandreamer: Fast generation from text to 3d gaussians by bridging 2d and 3d diffusion models. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 6796–6807. 2, 5, 6, 10

Yi, T.; Fang, J.; Zhou, Z.; Wang, J.; Wu, G.; Xie, L.; Zhang, X.; Liu, W.; Wang, X.; and Tian, Q. 2024b. GaussianDreamerPro: Text to Manipulable 3D Gaussians with Highly Enhanced Quality. *arXiv:2406.18462*. 2

Yu, X.; Guo, Y.-C.; Li, Y.; Liang, D.; Zhang, S.-H.; and Qi, X. 2023. Text-to-3D with Classifier Score Distillation. *arXiv.org*. 2

Zeng, B.-W.; Li, S.; Feng, Y.; Yang, L.; Li, H.; Gao, S.; Liu, J.; He, C.; Zhang, W.; Liu, J.; Zhang, B.; and Yan, S. 2023. IPDreamer: Appearance-Controllable 3D Object Generation with Complex Image Prompts. *arXiv preprint arXiv:2310.05375*. 2

Zhang, L.; Rao, A.; and Agrawala, M. 2023. Adding Conditional Control to Text-to-Image Diffusion Models. In *IEEE International Conference on Computer Vision (ICCV)*. 2, 5, 10

# AnchorDS: Anchoring Dynamic Sources for Semantically Consistent Text-to-3D Generation
# –Appendix–

## Implementation Details

**3D Representation.** We evaluate AnchorDS using both 3D Gaussian Splatting (3DGS) (Yi et al. 2024a) and NeRF as 3D representations to demonstrate robustness across different representation types. These representations exhibit different sensitivities to SDS guidance: 3DGS, being point-cloud-like, is highly sensitive to gradient flickering and often requires stabilization strategies like gradient clipping and hierarchical optimization. NeRF, while more stable, requires significantly more optimization steps to converge. For 3DGS experiments, we follow common practice by initializing with simple text-to-3D point clouds from Shap-E (Jun and Nichol 2023), as 3DGS cannot converge without reasonable initialization. This adds minimal computational cost due to the explicit nature of the representation. Importantly, unlike SDS-Bridge which requires SDS-guided rough initialization to prevent source estimate deviation, AnchorDS does not require such initialization. Our image-conditioned anchoring mechanism accurately captures the source distribution from the beginning, allowing direct optimization without preliminary SDS stages for NeRF experiments.

**AnchorDS Guidance Formula.** Adding back the variance reduction term $m_2$ in Eq. 6 in the main paper, our final AnchorDS guidance is:

$$\begin{aligned}
\mathcal{L}_{AnchorDS} &= g_t^{(\tau)} + m_2 \\
&= \hat{\boldsymbol{\epsilon}}_\phi(z_t; t, y) - \hat{\boldsymbol{\epsilon}}_\phi(z_t; t, \emptyset, I^{(\tau)}) \\
&\quad + \hat{\boldsymbol{\epsilon}}_\phi(z_t, t, \emptyset) - \boldsymbol{\epsilon}.
\end{aligned} \tag{15}$$

Negative prompts $y_{neg}$ may replace $\emptyset$ for more informative anchoring.

**Hardware & Training Setup.** All experiments are conducted on a single NVIDIA A40 GPU (48 GB). The threshold for the Filtering strategy is $\gamma = 0.03$, and the image adapter fine-tuning strategy adopts a learning rate of $1 \times 10^{-4}$. All remaining hyper-parameters mirror those of GaussianDreamer (Yi et al. 2024a) for the 3DGS pipeline and ProlificDreamer (Wang et al. 2023) for the NeRF pipeline.

To clarify the practical cost of our dynamic anchoring, we report wall-clock runtimes in Table 3. All runs are measured per text prompt. The additional image-conditioning pass in AnchorDS is executed in parallel with the original diffusion pass and thus incurs negligible overhead compared to SDS-based baselines. Optional filtering and lightweight adapter fine-tuning slightly increase runtime but remain a minor fraction of the total optimization.

## More Ablation Studies

To validate the robustness and generalizability of our approach across different prior model configurations, we conduct ablation studies examining the impact of various 2D diffusion models and image conditioning adapters. Quantitative results shown in Table 2 in the main paper also validate that AnchorDS is robust against the selection

Table 3: **Runtime comparison.** Wall-clock runtimes per text prompt on a single NVIDIA A40 GPU. AnchorDS matches the cost of SDS-based baselines; the optional filtering (Filter) and adapter fine-tuning (FT) introduce only marginal overhead while improving robustness and fidelity.

| Method | 3D Representation | Runtime / prompt |
|---|---|---|
| GaussianDreamer | 3DGS | 25 min |
| AnchorDS (ours) | 3DGS | 25 min |
| AnchorDS (ours) + Filter + FT | 3DGS | 30 min |
| ProlificDreamer | NeRF | 3.5 h |
| AnchorDS (ours) | NeRF | 3.5 h |
| AnchorDS (ours) + Filter + FT | NeRF | 4.0 h |

of different image conditionings, using both ControlNet with normal map and IP-Adapter with the identity image as image conditions achieve better results compared with the baseline. In Fig. 9, we show the qualitative evaluation of AnchorDS using SD 1.5 with three image conditioners—InstructPix2Pix (Brooks, Holynski, and Efros 2023) (IP2P), IP-Adapter (Ye et al. 2023), and ControlNet-NormalBae (Zhang, Rao, and Agrawala 2023). With SD 2.1, we adopt ControlNet-NormalBae conditioning.

**Cross-Model Performance:** Despite SD 1.5's more limited capabilities compared to SD 2.1, our method with IP-Adapter on SD 1.5 achieves competitive results with VSD running on the significantly more powerful SD 2.1 base model. When combined with SD 2.1 and ControlNet, our approach generates the most photorealistic textures among all evaluated configurations.

**Comparison with Existing Methods:** Across all tested configurations, AnchorDS variants consistently outperform existing methods including vanilla SDS, SDS-Bridge, and VSD, demonstrating the fundamental effectiveness of our dynamic source anchoring strategy regardless of the underlying architecture.

**Adapter Comparison:** Among our conditioning variations, both IP-Adapter and ControlNet produce high-quality and 3D-consistent results, confirming our architectural choice while validating AnchorDS's effectiveness across different conditioning mechanisms. IP-Adapter achieves superior performance in balancing the quality and 3D-consistency compared to other alternatives, Our method demonstrates strong generalizability across different architectural configurations, indicating significant potential for integration with future advanced diffusion models and emerging conditioning mechanisms as they become available.

## Discussion on Existing Methods under Our Formulation

Our approach offers a new perspective on source conditioning in text-to-3D generation. It is insightful to compare AnchorDS with recent techniques that also aim to address SDS's oversmoothing artifacts.

**Comparison with SDS-Bridge.** SDS-Bridge (McAllister et al. 2024) also recognizes the source distribution mismatch and proposes to use a negative prompt to describe the flaws of the current 3D model. In effect, SDS-Bridge replaces the
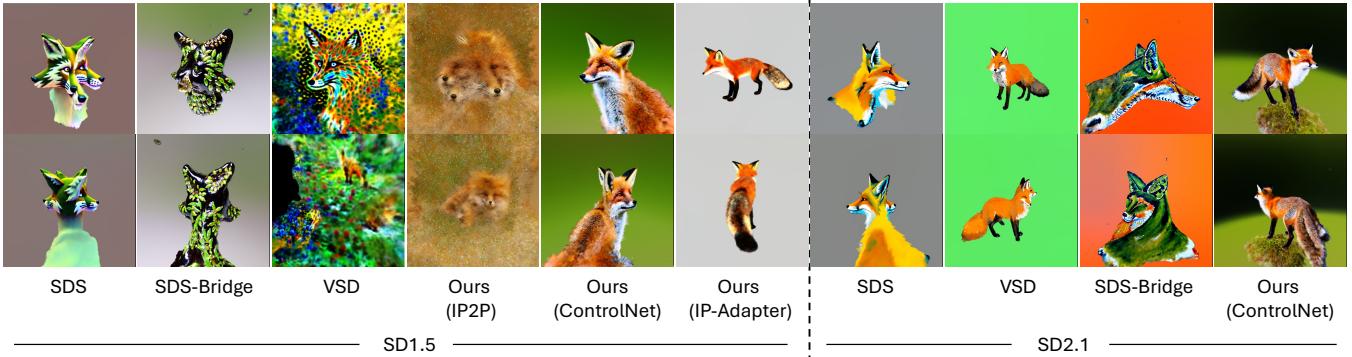
Figure 9: **Ablation studies across diffusion models and conditioning adapters.** Our AnchorDS variants consistently outperform existing methods (SDS, SDS-Bridge, VSD) across all configurations. Notably, AnchorDS with IP-Adapter on SD 1.5 achieves competitive quality with VSD on SD 2.1, while AnchorDS with ControlNet on SD 2.1 produces the most photorealistic results. Among conditioning adapters, IP-Adapter demonstrates superior 3D consistency compared to IP2P and ControlNet alternatives.

unconditional score with a negatively-conditioned one (*e.g*, a prompt describing "a bad, unfinished rendering"), then takes a difference similar to ours. However, this approach has fundamental limitations. A negative prompt is a *static* descriptor and may not accurately characterize the 3D state as it evolves. If the chosen negative prompt diverges from actual errors in renders, guidance can become misaligned, sometimes pushing results further off-track. Indeed, SDS-Bridge is typically applied only after a period of normal SDS optimization, to ensure the 3D model reaches a "rough" state that the negative prompt can plausibly describe. In contrast, AnchorDS uses the rendered image $I^{(\tau)}$ itself as the descriptor of the current state, which is by definition precise and up-to-date. By conditioning on $I^{(\tau)}$ at every iteration, our source estimate adjusts automatically as the 3D asset changes—treating the 3D generation as dynamic distribution alignment rather than one-shot static correction. This dynamic anchoring eliminates the need for separate SDS pre-runs or hand-crafted prompts describing the source; the model's render provides all necessary information. Empirically, we found AnchorDS to be more robust than SDS-Bridge, which can fail when negative prompts are inadequate.

**Comparison with ProlificDreamer.** Prolific-Dreamer (Wang et al. 2023) takes a different approach: rather than reweighting a pretrained model's guidance, it trains a specialized diffusion branch via LoRA finetuning to better represent the current 3D scene. Their VSD guidance is defined as the difference $\hat{\epsilon}_\phi(z_t; t, y_c) - \hat{\epsilon}_\phi^{(\text{LoRA})}(z_t; t, y, c)$, where $c$ represents view conditioning. While VSD also yields difference-of-noise guidance, the philosophy diverges significantly from ours. VSD fine-tunes a LoRA model to directly approximate the evolving source distribution in latent space, treating the source as a particle distribution of 3D parameters (implicitly aggregating multi-view images). Due to the dynamic nature of the source distribution, VSD's latent space approximation inherently lags behind the actual distribution. Each fine-tuning step provides only slow

incremental updates, resulting in persistent inaccuracies since the evolving distribution is never captured in time.

In contrast, AnchorDS does not aim to alter the diffusion model's internal latent distribution. Instead, it utilizes the diffusion model's existing capability to interpret and leverage image conditions. Fine-tuning an adapter in AnchorDS serves solely to familiarize the model with the conditional generation scenario rather than continuously updating internal distributions. Consequently, our fine-tuning is lightweight, converges quickly, and enables accurate, generalized estimation directly from the dynamically provided condition. This makes estimation inherently precise and immediately responsive to changes, avoiding the lag inherent in VSD's approach.

**Orthogonal Directions.** Several recent works address complementary aspects of SDS optimization. Methods like FSD (Yan, Wu, and Ma 2024) and SSD (Tang et al. 2024a) focus on variance reduction in the noise estimator $\epsilon$. ISM (Liang et al. 2024) addresses single-step score estimation inaccuracy through DDIM inversion for multi-step approximation. These approaches are orthogonal to our source distribution estimation focus and could potentially be combined with AnchorDS for further improvements.

## Comparison with Feed-forward and Hybrid Methods

**Qualitative Comparison with Trellis.** We qualitatively compare our method with Trellis (Xiang et al. 2024), a 3D latent-based text-to-3D model. As shown in Fig. 10, our results demonstrate superior geometric and color fidelity. In contrast, Trellis suffers from out-of-distribution generalization issues due to its reliance on 3D training datasets (He et al. 2023), and exhibits poor performance under text prompt conditions (Xiang et al. 2024).

**Quantitative Comparisons.** Beyond qualitative comparisons, we report CLIP-based text-image alignment scores on our 50-prompt benchmark in Table 4. We compare representative feed-forward text-to-3D methods against our An-

A black cat sleeps peacefully beside a carved pumpkin

A pair of white sneakers on a black mat

A bald eagle carved out of wood

A DSLR photo of a chow chow puppy

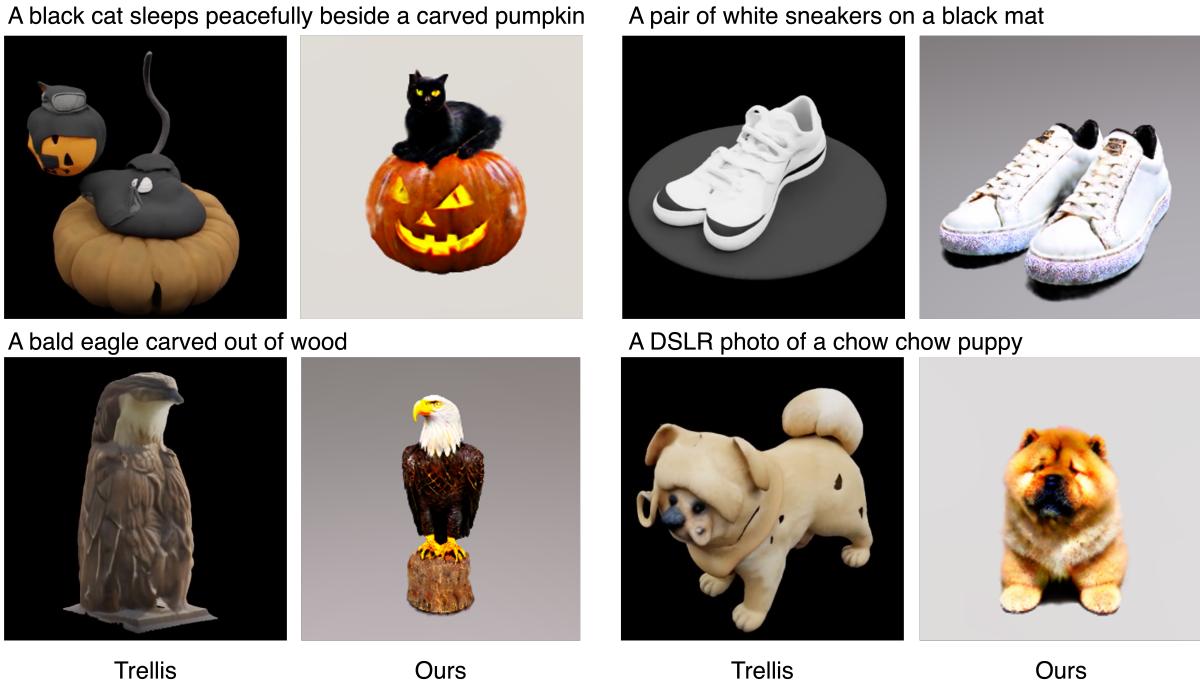Trellis          Ours          Trellis          Ours

Figure 10: **Qualitative comparison with Trellis.** Our results demonstrate superior geometric and color fidelity, while Trellis suffers from out-of-distribution generalization issues due to its reliance on specific training datasets (He et al. 2023), and exhibits poor performance under text prompt conditions (Xiang et al. 2024).

chorDS initialized from Shap-E (Jun and Nichol 2023). AnchorDS achieves substantially higher CLIP similarity than Trellis (Xiang et al. 2024), Shap-E, and Hunyuan3D-2.1, indicating stronger prompt fidelity despite operating as an SDS-based refinement. This supports our claim that dynamic source anchoring not only stabilizes SDS optimization but also surpasses existing feed-forward pipelines in semantic alignment.

Taken together, these quantitative results clarify the role of SDS-based optimization as a complementary paradigm to feed-forward pipelines. Feed-forward 3D generators are attractive for their fast inference, but are constrained by the coverage of their training corpora and often degrade on out-of-distribution or compositionally challenging prompts (He et al. 2023). In contrast, SDS operates directly on arbitrary 3D parameterizations and can be deployed as a retraining-free refinement module that post-improves feed-forward outputs. By stabilizing SDS through our state-anchored guidance, AnchorDS prevents drift and oversmoothing, turning SDS from a fragile heuristic into a robust mechanism for high-fidelity, controllable 3D generation.

## More Qualitative Results

Additional qualitative results for NeRF (Fig. 11) and 3DGS (Fig. 12) showcase consistent performance across diverse semantic categories and complexity levels.

Table 4: **Quantitative comparison with feed-forward and hybrid methods.** CLIP-based image-text similarity on our 50-prompt set. AnchorDS consistently outperforms representative feed-forward and hybrid approaches, demonstrating that dynamic source anchoring yields superior semantic alignment while retaining the flexibility of SDS-style optimization.

| Method | CLIP image-text sim ↑ |
|---|---|
| Shap-E (Jun and Nichol 2023) | 0.25 |
| Trellis (Xiang et al. 2024) | 0.22 |
| Hunyuan3D-2.1 (Team 2025) | 0.29 |
| AnchorDS (ours) | **0.37** |

## Details of User Study

We conducted a human preference study to evaluate the effectiveness of our proposed AnchorDS compared to existing methods using all 50 complex prompts. A total of 912 unique participants were recruited through Amazon Mechanical Turk, resulting in 1000 effective comparison samples in total. The interface is shown in Fig. 13. For each comparison, participants were provided with the text prompt and the randomly ordered generated results of various methods. They were then asked to indicate their preference rankings based on three evaluation criteria:

- 3D Consistency: The output that maintains the best 3D consistency;
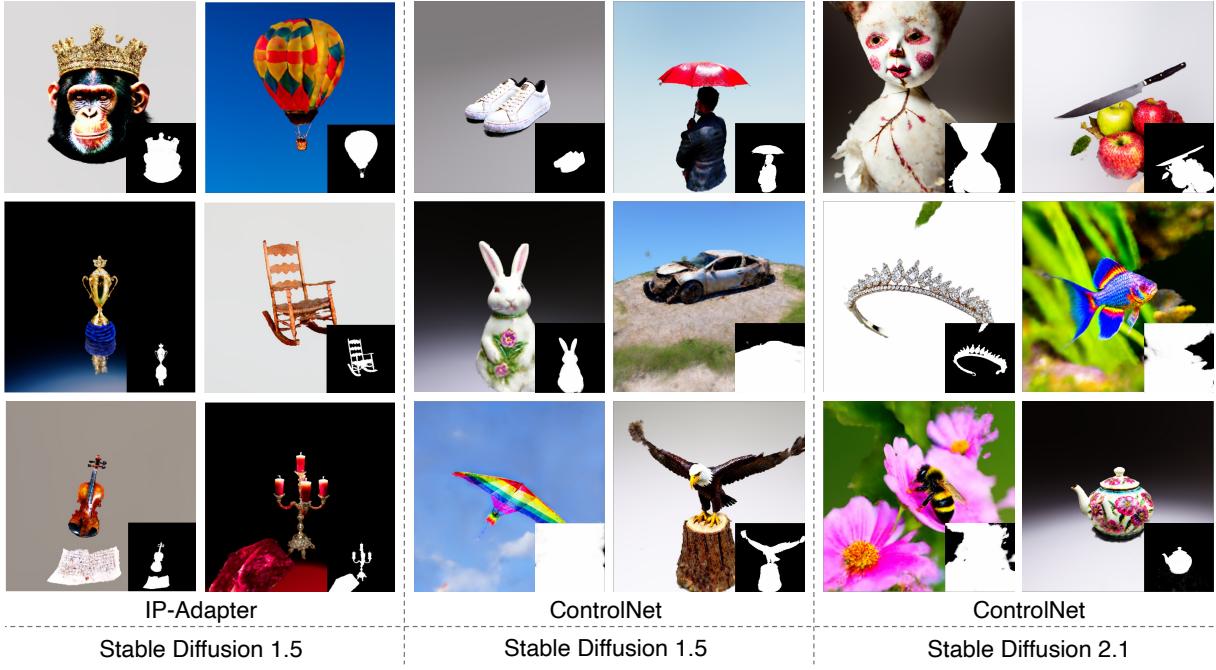- Test Consistency: The output that better maintains consistency with the input prompt;

Figure 11: **More qualitative results on text-to-NeRF.** Prompts include: "A chimpanzee dressed like Henry VIII king of England", "A hot air balloon in a clear sky", "A pair of white sneakers on a black mat", "A man is holding an umbrella against rain", "A cracked porcelain doll's face", "Ripe apples cluster next to a gleaming knife", "The golden trophy shines brightly next to a ruffled blue ribbon", "A wooden rocking chair on a porch", "A small porcelain white rabbit figurine", "A completely destroyed car", "A sparkling diamond tiara", "A beautiful rainbow fish", "A violin reclines on a chair next to a music sheet filled with notes", "A DSLR photo of a candelabra with many candles on a red velvet tablecloth", "A rainbow-colored kite soaring in the sky", "A bald eagle carved out of wood", "A bumblebee sitting on a pink flower", "A ceramic teapot with floral patterns".

- Visual Quality: The output with the best overall visual quality.

The methods are grouped as SD 2.1-based and SD 1.5-based, and evaluated separately. For the SD 2.1-based methods, the users are simply asked to indicate their most preferred result, since there are just two methods to compare.

## Limitations and Future Work

Our method shares several limitations common to SDS-based approaches. In particular, the quality of the generated 3D content is highly dependent on the guidance from 2D prior models; if the underlying 2D model fails to generate meaningful representations, the corresponding 3D outputs are also compromised. This limitation could be alleviated by leveraging more powerful 2D backbones, such as SD3 or Flux.1. Additionally, we observe that convergence stability is sensitive to the choice of 3D representation, suggesting room for improvement in representation learning and rendering fidelity.

On the score distillation side, our work focuses primarily on addressing the first type of error—source estimation bias—as discussed in SDS-Bridge. We leave the integration of methods to mitigate the second class of error—trajectory estimation mismatch—as future work. Promising directions include the use of inversion-based techniques or improved sampling strategies to more accurately track the evolving distribution across denoising steps.

Finally, our key contribution lies in reinterpreting 3D generation as an evolving editing process. This perspective opens up new avenues for future research, such as unifying the formulation of 3D generation and 3D editing within a single framework, and extending our approach to more controllable or user-driven editing tasks.

Figure 12: **More qualitative results on text-to-3DGS.** The figure showcases 3D object generation results using various prompts: "A green cactus in a clay pot", "A red barn in a green field", "A lighthouse on a rocky shore", "A vibrant orange pumpkin sitting on a hay bale", "A vintage clock hanging on a brick wall", "A yellow school bus on a city street", "A bright red fire hydrant", "A cactus with pink flowers", "A gold glittery carnival mask", "A neon green skateboard with black wheels", "A well-worn straw sun hat", "A pair of shiny black patent leather shoes", "A pirate flag with skull and crossbones", "A vibrant, handmade patchwork quilt", and "Hot popcorn jump out from the red striped popcorn maker". Each result demonstrates the method's ability to generate diverse 3D objects with varying complexity, materials, and semantic categories.



Figure 13: Snapshot of the user study interface. Participants were shown the prompt and rendered images generated by different methods. They were asked to rank the methods following their preferences based on overall quality, 3D consistency, and consistency with the prompt.