## Turbo3D: Ultra-fast Text-to-3D Generation

<sup>1</sup> Carnegie Mellon University <sup>2</sup> Massachusetts Institute of Technology <sup>3</sup> Adobe Research

https://turbo-3d.github.io/

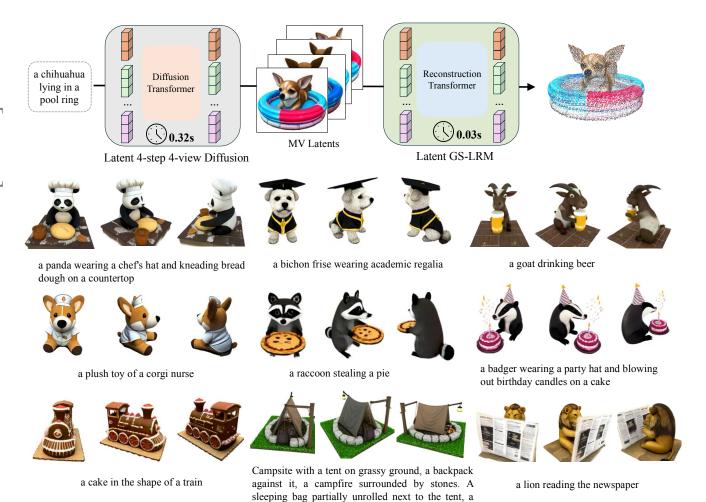


Figure 1. **Overview of our Turbo3D text-to-3D system.** Turbo3D generates high-quality 3D Gaussian Splatting (3DGS) assets from user prompts in *less than 1 second* on a single A100 GPU. It's a two-stage pipeline consisting of a highly efficient latent-space few-step multi-view (MV) generator and single-step MV reconstructor. Note that we visualize latents as RGB images and 3DGS assets as point clouds in the pipeline figure for clarity.

lantern hanging from a tree branch

<sup>\*</sup>Equal advising.

#### **Abstract**

We present Turbo3D, an ultra-fast text-to-3D system capable of generating high-quality Gaussian splatting assets in under one second. Turbo3D employs a rapid 4-step, 4-view diffusion generator and an efficient feed-forward Gaussian reconstructor, both operating in latent space, as shown in Fig. 1. The 4-step, 4-view generator is a student model distilled through a novel Dual-Teacher approach, which encourages the student to learn view consistency from a multi-view teacher and photo-realism from a single-view teacher. By shifting the Gaussian reconstructor's inputs from pixel space to latent space, we eliminate the extra image decoding time and halve the transformer sequence length for maximum efficiency. Our method demonstrates superior 3D generation results compared to previous baselines, while operating in a fraction of their runtime.

#### 1. Introduction

The recent advances in image generative models allow users to generate detailed outputs from just a text prompt. While initial denoising diffusion-based methods [7, 33, 35, 36, 43, 49] enabled impressive photo-realistic generation, recent techniques [18, 21, 22, 38, 44, 59, 60] have significantly improved inference efficiency of such models, allowing high-fidelity generation in the blink of an eye. Unfortunately, these advances in generative 2D modeling methods have not yet been matched in the 3D domain, where ultra-fast realistic 3D generation remains a challenge. In this work, we seek to bridge this gap, and present Turbo3D, a text-to-3D generative model that can synthesize detailed 3D outputs in a fraction of a second.

The existing approaches for text-based 3D inference can be categorized as either generative [9, 16, 23, 29, 56, 62] or optimization-driven [17, 31, 54]. The latter class of methods optimize 3D representations by 'distilling' pretrained 2D diffusion models [31]. While this approach can yield decent 3D outputs, it is highly inefficient, typically requiring several minutes or even hours to output a single 3D representation. The alternative paradigm is to learn a generative model that directly outputs 3D representations. While initial methods [9, 27] investigated representations such as point clouds and SDFs, the limited availability of 3D data restricted the generation quality. Recent approaches [16, 23, 39, 40] have instead advocated for learning generative models of multi-view images (followed by deterministic 3D reconstruction), as these can be initialized from pre-trained 2D generative models. While these methods have resulted in impressive generations, the multiview finetuning on synthetic data does inhibit their generation quality. More crucially, their inference efficiency is restricted by the iterative denoising process required for the text-conditioned multi-view generation.

In this work, we follow the paradigm of text-to-3D via multi-view generation, and aim to improve the efficiency of the underlying components to enable ultra-fast generation, while enhancing the fidelity of the generated outputs. Inspired by the recent progress in reducing inference time by distilling 2D diffusion models into one-step or few-step generators [37, 38, 44, 60], we adapt these for multi-view 3D generation. We first train a many-step text-to-multiview diffusion model and then distill it into a much faster 4-step generator using the distribution matching distillation (DMD) loss [60]. However, we find that this process leads to a significant quality degradation, as it fails to capture the full range of modes present in the multi-view teacher model. To overcome this, we propose to extend the DMD pipeline to incorporate another single-view teacher – a 2D denoising diffusion model trained on a large set of high quality aesthetic images. Our few-step multi-view generator is thus trained with a dual-teacher distillation approach, where the multi-view DMD loss helps our model learn multi-view consistency, and the single-view DMD loss ensures highfidelity outputs. To further improve the 3D generation efficiency, we note that our multi-view generator outputs latent representations (and not pixels) for the multi-view images. We build on this insight to adapt a prior multi-view to 3D reconstruction approach [61] to instead consume multi-view latents as input, and show that this improves the reconstruction efficiency without any performance loss.

Our overall system thus combines an efficient few-step multi-view generator with a multi-view latent to 3D model to enable ultra fast 3D generation. We train our model on the subset of Objaverse dataset [3, 4], which contains about 400k instances with Cap3D text captions [24]. We show that our proposed Turbo3D is able to produce high-quality 3D assets in less than one second (see Fig. 1), and also achieves comparable quality with previous state-of-the-art.

#### 2. Related Work

In this section, we discuss the closely related prior feed-forward 3D generation methods. Optimization-based methods [17, 31] is out of this work's scope. We review feed-forward 3D generative models in the following three categories: 1) methods that directly generate full 3D representation encoding geometry and appearance; 2) methods that generate shapes and then generate the textures; 3) methods that generate multi-view (MV) images followed by reconstruction. We also review the diffusion distillation literature which our method is built upon.

## 2.1. Directly Genereating Full 3D Representation

Several prior methods [1, 9, 29, 45, 47, 53, 57] have been introduced that directly generate 3D representations encoding both geometry and texture, e.g., implicit fields [9], point

clouds [29], and triplanes [53]. This line of work typically have to preprocess the source 3D data (meshes or multiview images) into the target representations for generative models in a lossy fashion, hence limiting their quality and scalabilty.

Alternative methods [1, 45, 47, 57] have been proposed to get rid of the lossy preprocessing. They implicitly bake a 3D generative process into a multi-view diffusion framework, creating a single-stage 3D generative model. However, trained only on 3D data, they suffer from the issue of limited generalization, lower pixel quality and reduced understanding of complex text prompts, compared with methods leveraging powerful pretrained text-to-image models [16, 20].

Our method leverages the strong priors in a pretrained image generative model, but we go a step further by distilling the slow teacher model into a fast student generator for ultra-fast text-to-3D.

## 2.2. Shape Generation + Texture Generation

Recent advancements in 3D content creation have introduced innovative methods [41, 56, 62] that focus on generating high-quality 3D shapes first, followed by generative texturing [2, 34], rather than creating both simultaneously. CLAY [62] offers a framework using a multiresolution Variational Autoencoder (VAE) and a latent Diffusion Transformer (DiT) to initially create detailed 3D geometries from inputs like text and images, before applying high-resolution physically-based rendering (PBR) textures. Direct3D [56] enhances scalability in image-to-3D generation by improving the performance of 3D VAE and DiT for generating 3D shapes, subsequently adding textures.

However, besides the challenge of learning a compact latent space suitable for generation, these methods are also usually slow as a result of the two diffusion models - one for shapes and the other for textures. In contrast, our Turbo3D only have one diffusion model for generating multi-views and is fast once distilled, while the reconstructor is deterministic and hence fast.

#### 2.3. MV Generation + MV Reconstruction

A large body of the 3D generation work [16, 19, 23, 40] focus on leveraging multi-view generation and reconstruction to enhance quality and efficiency. Instant3D [16] offers a fast method for creating high-quality 3D assets from text prompts by combining sparse-view generation with a transformer-based reconstructor [8, 52, 55, 61], significantly reducing inference time. MVDream [40] employs a multi-view diffusion model to improve consistency and stability in 3D generation, integrating 2D and 3D data. SyncDreamer [23] enhances multiview-consistency from a single-view image using a 3D-aware feature attention mechanism. One-2-3-45 [19] introduces an efficient ap-

proach to single image 3D reconstruction without extensive optimization, producing consistent 3D meshes. SV3D [50] utilizes a latent video diffusion model for novel multi-view synthesis, incorporating explicit camera control to improve 3D reconstruction quality.

We follow the same approach of reconstructing generated multi-views to create 3D assets in this work, but we focus on improving the efficiency of such systems. With our novel Dual-Teacher Distillation and Latent GS-LRM components, our Turbo3D manages to be at least an order of magnitude faster than these baselines while maintaining competitive quality. Concurrent work GECO [51] also use diffusion distillation to speed up image-to-3D. However, we differ in our dual-teacher distillation design with focus on text-to-3D. Our pipeline is also simpler and avoids the cumbersome mesh reconstructions for 3D distillation in GECO.

#### 2.4. Diffusion Distillation

Recent progress in diffusion models have focused on improving the efficiency of image generation by reducing the number of sampling steps required, leading to the development of several innovative distillation techniques [25, 26, 28, 37, 38, 44, 58-60]. Methods like Improved Rectified Flows [14] and InstaFlow [22] straighten the ODE trajectories, making them easier to approximate with a one-step student model. Consistency Models [44] train student generators to map any point on the teacher's ODE trajectory to a consistent target, enabling one-step and few-step image generation. Distribution Matching Distillation (DMD) [60] trains a one-step generator by minimizing the reverse KL divergence between the data distribution and the generator's output distribution [5, 25, 54]. Building on this, DMD2 [59] integrates GAN losses [6] and extends the method to multistep generators, further enhancing generation quality.

However, most of existing approaches focus on distilling pretrained 2D image diffusion models. In our work, we adopt the popular DMD [60] approach and extend it to the multi-view domain. Distilling multi-view models presents unique challenges, such as increased mode collapse due to fine-tuning and distillation processes. To address this issue, we introduce a novel dual-teacher distillation technique, enabling swift, photorealistic multi-view generation.

## 3. Background

#### 3.1. Multi-view Diffusion Model

Diffusion models generate data from noise by reversing a forward noising process. Given a data sample  $x_0 \sim p(x_0)$  from the data distribution, the forward process progressively adds Gaussian noise over T timesteps. At timestep t, the noised distribution conditioned on  $x_0$  is given by:

$$q(x_t|x_0) = \mathcal{N}(x_t; \alpha_t x_0, \beta_t^2 \mathbf{I}), \tag{1}$$

where  $\alpha_t$  and  $\beta_t$  are timestep-dependent constant that are specified by the noise schedule [11, 43]. The diffusion model  $\epsilon_{\theta}$  is trained to reverse this noising process by learning to predict the added Gaussian noise  $\epsilon$ :

$$L(\theta) = \mathbb{E}||\epsilon - \epsilon_{\theta}(x_t, t)||^2.$$
 (2)

There exists other formulations, in which the diffusion model learns to predict the clean input  $x_0$  directly [10], or a combination of  $x_0$  and  $\epsilon$  [37]. Regardless of the specific prediction target, the outputs of these models can be related to the score function, which is the gradient of the log probability density of the data distribution [43]:

$$s_{\theta}(x_t, t) = \nabla_{x_t} \log p(x_t) = -\frac{\epsilon_{\theta}(x_t, t)}{\beta_t} = -\frac{x_t - \alpha_t x_{\theta}(x_t, t)}{\beta_t^2}$$
(3)

where  $\epsilon_{\theta}$  [7] and  $x_{\theta}$  [10] represents noise and data prediction diffusion models, respectively. In our paper, we adopt the noise prediction scheme but our method generalizes to any formulations.

In multi-view diffusion models, the generation of a 3D scene is conditioned on a text prompt, enabling the joint denoising of multiple views to produce a set of 3D-consistent output images [16, 40]. Specifically, given a text prompt c and a set of multi-view images  $\{x^i\}_{i=1}^K$  where K is the number of view, the multi-view diffusion model learns to predict the added noise across all views simultaneously. The training objective is formulated as:

$$\mathbb{E}||\epsilon - \epsilon_{\theta}(\{x_t^i\}_{i=1}^K, t, c)||^2. \tag{4}$$

where  $\epsilon$  represents independent Gaussian noise with the same variance applied to each view. At inference time, the generation starts from a set of fully noisy multi-view images  $\{x_T^i\}_{i=1}^K$  sampled from a standard Gaussian distribution. The multi-view diffusion model iteratively generate a sequence of cleaner multi-view images. Various diffusion samplers [7, 42] can be employed during this process to optimize generation speed and quality.

## 3.2. Distribution Matching Distillation

Distribution Matching Distillation (DMD) is a widely used diffusion distillation technique that converts teacher diffusion models into a student generator requiring significantly fewer sampling steps [59, 60]. The DMD approach trains the student generator  $G_{\theta}$  by minimizing an approximate reverse KL divergence between the smoothed student's output distribution (denoted as  $p_{\text{fake}}$ ) and the smoothed data distribution (denoted as  $p_{\text{real}}$ ):

$$L_{\text{DMD}}(\theta) = D_{\text{KL}}(p_{\text{fake}}||p_{\text{real}}) = \mathbb{E}_{x,t} \left( \log(\frac{p_{\text{fake}}(x_t)}{p_{\text{real}}(x_t)}) \right). \quad (5)$$

Let the score functions of the data distribution and the student's output distribution be denoted as  $s_{\rm real}$  and  $s_{\rm fake}$ , respectively. The gradient of this KL divergence can be effectively approximated by the difference between these two

score functions:

$$\nabla_{\theta} L_{\text{DMD}}(\theta) \approx \mathbb{E} \left[ -\int \left( s_{\text{real}} \left( F(G_{\theta}(\epsilon), t), t \right) - s_{\text{fake}} \left( F(G_{\theta}(\epsilon), t), t \right) \right) \frac{dG_{\theta}(\epsilon)}{d\theta} d\epsilon \right]$$
(6)

where F represents the forward diffusion process defined in Eq. 1. The student generator can be adapted for a multi-step generation setting by replacing the pure noise input  $\epsilon$  with a partially noisy image  $x_t$  [59]. During training, the data distribution's score function is initialized from the teacher diffusion model and kept fixed, while the student's output distribution score function is dynamically trained using the student's output and a denoising loss (Eq. 4).

#### 4. Method

In this section, we provide the technical details of our Turbo3D text-to-3D system, which features a highly efficient multi-view (MV) generator and reconstructor. We begin by describing our novel DualTeacher Distillation approach; it creates a rapid MV generator by jointly distilling knowledge from both a multi-view (MV) teacher and a single-view (SV) teacher. Following this, we discuss the latent-space GS-LRM that instantly lifts the generated MV latents into high-quality 3D Gaussians.

#### 4.1. Dual-teacher Distillation for MV Diffusion

Leveraging a multi-step multi-view diffusion in a text-to-3D generation process can be inefficient due to repeated evaluations of the diffusion denoiser in the sampling process. To speed this up, one approach is to use diffusion distillation [59, 60] to train a single-step or few-step MV generator.

However, we find that naively applying diffusion distillation methods to a MV teacher can cause the student model to generate overly simplistic and cartoonish appearance, which closely resembles the 3D Objaverse dataset used during MV teacher finetuning and distillation [3] (shown in Fig. 6). We call this phenomenon 'compounding mode collapse'; this happens because both finetuning and distillation sacrifices generation diversity for efficiency. As the MV teacher is already biased towards synthetic Objversestyle appearance, further distilling it will have compound effect of locking the distilled generator in the mode of Objaverse data that are far from the modes of photorealistic natural images.

To address this issue, we propose to use dual teachers in the distillation process: one MV teacher to teach the student model about multi-view consistency, and one SV teacher to teach about each views' photo-realism. We illustrate this Dual-teacher Distillation algorithm in Fig. 2. Concretely, this is formulated as:

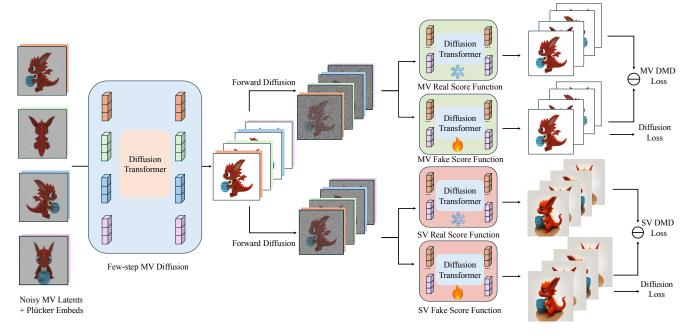


Figure 2. **Dual-teacher distillation framework in our Turbo3D.** Note that latents are visualized as RGB images for clarity. We aim to distill a multi-step multi-view teacher generator (right, green) into a few-step multi-view generator (left, blue). Our few-step MV student generator is conditioned on Plücker embeddings for better 3D awareness. Similar to [59], we optimize the student generator using distribution matching objective (DMD loss) and train the fake score function to model the distribution of samples produced by the student generator. In particular, we integrate two teacher models: multi-view teacher and single-view (SV) teacher to enhance both multi-view consistency and photorealism. The MV score functions take a set of images of one object as input and calculate the MV DMD loss, while the SV score functions treat each image separately and calculate the SV DMD loss.

$$L_{\text{DMD}}^{\text{Dual}}(\theta) = D_{\text{KL}}\left(p_{\text{fake}}(\{x_t^i\}_{i=1}^K) \mid\mid p_{\text{real}}^{\text{MV}}(\{x_t^i\}_{i=1}^K)\right) + \lambda \cdot \frac{1}{K} \sum_{i=1}^K D_{\text{KL}}\left(p_{\text{fake}}(x_t^i) \mid\mid p_{\text{real}}^{\text{SV}}(x_t^i)\right), \quad (7)$$

where  $p_{\rm real}^{\rm MV}(x_t), p_{\rm real}^{\rm SV}(x_t)$  represent MV and SV teachers, respectively;  $\lambda$  is the loss weight balancing the influence of MV and SV teachers on the distilled student model; K=4 is the number of views. We set  $\lambda$  to 1 in our experiments. As demonstrated in Fig. 6, having the additional SV teacher in the distillation process effectively addresses the compound mode collapse issue, because it tries to pull each view to look like natural images.

#### 4.2. Latent GS-LRM for MV Reconstruction

To reconstruct 3D from the generated MV latents, one straightforward approach is to decode them into multi-view images, and then use the pixel-space GS-LRM [61] to produce the 3D Gaussians. However, such an approach can suffer from efficiency and memory issue when scaling to high resolution due to the poor performance of Conv2D operators in VAE decoder [12].



Figure 3. We compare the renderings of pixel GS-LRM and latent GS-LRM. Latent GS-LRM achieves comparable reconstruction quality as pixel GS-LRM.

We propose to skip the VAE decoding and directly input the generated MV latents to a latent GS-LRM for best efficiency. To train such a model, we supervise the reconstructed Gaussians with pixel-space novel-view rendering losses ( $\ell_2$  and perceptual losses as in [61]). We show that replacing pixel-space GS-LRM with a latent one does not affect the quality of generated assets in Tab. 3 and Fig. 3, while being faster.

## 5. Experiments

In this section, we first describe our experimental setup (Section 5.1). We then compare our method with state-of-the-art text-to-3D baselines (Section 5.2). Finally, we ablate each component of our framework to showcase their effectiveness (Section 5.3).

## **5.1. Experimental Setup**

**Datasets.** We use the Objaverse dataset [3] to train both our multi-view generation model and multi-view reconstruction model. We scale the objects and center them to fit into a cube  $[-1,1]^3$ . For the generation task, we render the dataset at a fixed elevation (20 degrees) and 16 equidistant azimuths to achieve a good coverage of the object. We render using a field of view  $50^\circ$  at a distance of 2.7 and uniform lighting. For the reconstruction task, we render 32 views randomly placed around the object with a random distance in the range of [1.5, 2.8]. We render a total of 730K objects. **Baselines** We adopt Instant3D [16] and I GM [46] as

**Baselines.** We adopt Instant3D [16] and LGM [46] as our baseline text-to-3D methods. However, since the field of fast text-to-3D is relatively underexplored, we also include recent state-of-the-art fast image-to-3D methods TripoSR [48] and SV3D [50] as baseline methods. For a fair comparison, we use a popular few-step text-to-image model Flux [13] to generate an input image first and then feed it to the image-to-3D models. The inference time of image-to-3D models is a summation of the two parts.

Metrics. We adopt the CLIP score [32] and VQA score [15] to assess the semantic alignment between the generated results and text prompts. We use 400 text prompts from DreamFusion [31] for evaluation. We generate one object for each prompt. For each generated 3D object, we render 10 random views and calculate the average CLIP score and VQA score between the rendered images and the input text. For inference time, we report it with all methods under the same image resolution of 256. Notably, some methods, e.g., Instant3D [16], only support a higher resolution of 512. For a fair comparison, we report their quantitative results under the resolution of 512 and only inference time on the resolution of 256.

Implementation Details. The whole pipeline of our method includes three training phases. We first train a multi-step multi-view diffusion model on the Objaverse dataset, by fine-tuning an internal DiT [30] based text-to-image model. We train this model for 30k iterations with 32 80G A100 GPUs using a total batch size of 128 and a learning rate of  $3e^{-5}$ . Then we perform distillation to distill the multi-step multi-view generator into a few-step multi-view generator, which takes 10k iterations with 32 80G A100 GPUs with a global batch size of 128 and a learning rate of  $5e^{-6}$ . Finally, we train a reconstruction model – latent GS-LRM – from scratch, which takes 80k iterations with 32 80G A100 GPUs with a total batch size of 256 and a

Method	Clip	VQA	Inference
	Score ↑	Score ↑	Time ↓
TripoSR [48]	23.85	0.57	1.19s
SV3D [50]	24.92	0.64	12.52s
Instant3D [16]	26.23	0.65	15.02s
LGM [46]	24.73	0.58	6.56s
Turbo3D (Ours)	<b>27.61</b>	<b>0.76</b>	<b>0.35</b> s

Table 1. Comparison against state-of-the-art 3D generation methods. Our Turbo3D generates 3D assets with highest CLIP and VQA scores while using the least amount of time (benchmarked on a A100 GPU).

learning rate of  $4e^{-4}$ .

## 5.2. Evaluation against baselines

**Oualitative Comparisons.** As shown in Fig. 4, our method generates significantly better results compared with LGM [46] and Instant3D [16]. In particular, LGM tends to produce simple 3D assets without geometry and texture details; it also lacks robustness, oftentimes generating low-quality 3D assets that are broken, or suffer from multi-face Janus problem [31], or do not closely follow provided text prompts. Instant3D is much more robust and able to produce plausible 3D assets most of the time. However, its performance on text-image coherence is limited compared with our Turbo3D. For example, in the first example, LGM's rendition appears oversimplified with indistinct egg and yolk boundaries, while Instant3D provides a closer match but lacks fine details in the egg's structure and wood texture and fails to capture the concept of 'spill out'. For complex prompts that describes a scene with multiple objects (4th and 6th rows), LGM generates a Janus asset while Instant3D miss a lot of concepts like blanket, basket, wildflowers etc; our generations adhere to the prompts much more closely. These qualitative results highlight our Turbo3D's ability to generate high-quality, text-aligned 3D models with a superior level of detail, realism, and coherence across diverse and complex prompts, outperforming both LGM and Instant3D by a large margin.

Quantitative Comparisons. Tab. 1 presents a quantitative comparison of our proposed method, Turbo3D, against several state-of-the-art approaches, including TripoSR [48], SV3D [50], Instant3D [16], and LGM [46], on the text-to-3D task. Our proposed Turbo3D achieves the highest CLIP Score of 27.61 and VQA Score of 0.76, outperforming other methods by a significant margin in both quality metrics. In addition to the quality improvement, Turbo3D demonstrates remarkable efficiency with an inference time of only 0.35 seconds, substantially faster than competing methods. Although TripoSR is able to generate a 3D asset in only 1.19s, the quality of the generated results is highly limited. As a result, our proposed Turbo3D is able to achieve outstanding



Figure 4. Comparison of our Turbo3D against baselines LGM [46] and Instant3D [16]. Among these methods, Our method generates the most detailed and physically plausible 3D assets, closely adhering to the provided text prompts. In contrast, LGM tends to generate broken assets with Janus issue [31], while Instant3D has poorer text alignment, oftentimes missing some concepts, e.g., 'spilling out' in the first row, 'river' in the second row, etc.

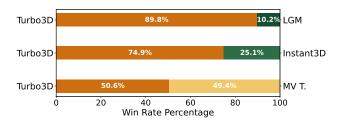


Figure 5. User study results comparing our Turbo3D to baseline LGM [46], Instant3D [16], and our slow MV teacher. Our Turbo3D is consistently preferred over baseline LGM and Instant3D, while having on-par preference with our MV teacher. See Fig. 4,6 for visual comparison.

performance in terms of both quality and inference speed compared with state-of-the-art methods.

**User Study.** We run a user study by randomly selecting 80 text prompts and asking 56 users to make 1120 pairwise comparisons (users are shown an input text prompt and two

generated 3D assets from two anonymous methods). We show results in Fig. 5. Since the quantitative results of TripoSR is not competitive with others, it is not included in the user study. Moreover, SV3D only outputs videos instead of 3D representations of objects, so it is hard to perform a fair comparison. Therefore, we compare our Turbo3D with LGM [46], Instant3D [16], and our MV Teacher. In particular, MV Teacher is the model which our Turbo3D gets distilled from. When compared to LGM, Turbo3D achieved a win rate of 89.8%, with only 10.2% of participants favoring LGM. Against Instant3D, Turbo3D also outperformed with a 74.9% win rate, indicating that users consistently found Turbo3D's outputs more aligned with the input text descriptions and more visually compelling. When evaluated against our teacher model, Turbo3D held a close win rate of 50.6%, with 49.4% preferring MV T., reflecting comparable quality between the student and teacher models. The re-



MV teacher (10.18s) MV teacher only distillation (0.35s) Dual teacher distillation (0.35s)

Figure 6. Ablation of our Dual-teacher distillation algorithm. Naively distilling MV teacher (middle column) causes compound mode collapse (see Sec. 4.1), producing overly smooth synthetic-looking assets. Our dual-teacher distillation (right column) fixes the issue and generates 3D assets that are as photorealistic as, if not more than, the baseline MV teacher (left column). We also include the inference timings for each method; the distilled model is  $\sim 50x$  faster than the teacher model.

sult indicates that our Turbo3D not only achieves significant speedup with distillation, but also preserves the generation ability from the teacher model.

#### **5.3. Ablation Study**

**Effect of Single-view Teacher.** In Tab. 2, we demonstrate the effectiveness of the dual teacher distillation strategy. The first line is our multi-view (MV) teacher model, which can achieve impressive results but runs slowly because of the many diffusion sampling steps required. When performing distillation only with this MV Teacher, the quality drops by a large margin as shown in the second row. When adding a single-view teacher model for distillation, the distilled model is able to achieve much better results compared with the previous one. This configuration approaches the performance of the Multi-step MV Model while maintaining the efficiency benefits of the few-step setup, showcasing the advantages of using a dual-teacher strategy in our distillation. We also provide visualizations of the three models in Fig. 6. These comparisons demonstrate that the dual teacher distillation model strikes a balance between detail retention and stylization, closely replicating the quality of the MV teacher model while benefiting from the efficiency gains of distillation.

**Effect of Latent GS-LRM.** In Tab. 3, we showcase the effectiveness of latent GS-LRM. Compared with the original GS-LRM which operates in pixel space, our latent GS-LRM is able to skip the expensive image decoding process while achieving similar image quality.

Ablation	CLIP Score ↑	VQA Score ↑
Multi-step MV Model	28.04	0.77
Few-step Model (MV Teacher) Few-step Model (Dual Teacher)	26.60 27.61	0.69 0.76

Table 2. Ablation of dual teacher distillation. Distillation leads to quality drop compared with the MV teacher model. Compared with naively distilling the single MV teacher, dual-teacher distillation leads to much smaller quality drop. See Fig. 6 for visual comparison.

Ablation	CLIP	VQA	Inference
	Score ↑	Score ↑	Time ↓
Pixel GS-LRM [61]	<b>27.62</b> 27.61	0.76	0.45s
Latent GS-LRM		0.76	<b>0.35s</b>

Table 3. **Ablation for latent GS-LRM.** We report the CLIP score, VQA score, and overall text-to-3D inference time for comparison. Our latent GS-LRM achieves similar image quality while enabling better efficiency ( $\sim$ 22% speedup).

#### 6. Conclusion

In this work, we propose Turbo3D for ultra-fast text-to-3D generation. To enable fast multi-view generation, we pro-

pose to distill a multi-step multi-view generator into a fewstep multi-view generator. Moreover, to restore the multiview consistency and photo-realism during distillation, we introduce a novel dual-teacher distillation framework. To further improve the multi-view reconstruction efficiency, we propose a latent GS-LRM which directly reconstructs 3D Gaussians from multi-view latents. Extensive experiments demonstrate that our proposed Turbo3D is able to achieve outstanding performance in terms of both generation quality and inference efficiency.

## Acknowledgements

This work began during Hanzhe Hu and Tianwei Yin's internships at Adobe Research. This research was supported in part by NSF award IIS-2345610.

## References

- [1] Titas Anciukevičius, Zexiang Xu, Matthew Fisher, Paul Henderson, Hakan Bilen, Niloy J Mitra, and Paul Guerrero. Renderdiffusion: Image diffusion for 3d reconstruction, inpainting and generation. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 12608–12618, 2023. 2, 3
- [2] Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner. Text2tex: Text-driven texture synthesis via diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18558–18568, 2023. 3
- [3] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. 2, 4, 6, 12
- [4] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of 10m+ 3d objects. Advances in Neural Information Processing Systems, 36, 2024. 2
- [5] Jean-Yves Franceschi, Mike Gartrell, Ludovic Dos Santos, Thibaut Issenhuth, Emmanuel de Bézenac, Mickaël Chen, and Alain Rakotomamonjy. Unifying gans and score-based diffusion as generative particle models. In *NeurIPS*, 2023. 3
- [6] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. Advances in neural information processing systems, 27, 2014. 3
- [7] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020. 2, 4
- [8] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint arXiv:2311.04400*, 2023. 3

- [9] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. arXiv preprint arXiv:2305.02463, 2023.
- [10] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *NeurIPS*, 2022. 4
- [11] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. In *NeurIPS*, 2021. 4
- [12] MIT HAN Lab. Patch conv: Patch convolution to avoid large gpu memory usage. https://hanlab.mit.edu/blog/patch-conv, 2024. Accessed: 2024-11-14. 5
- [13] Black Forest Labs. Flux. https://github.com/ black-forest-labs/flux, 2024. Accessed: 2024-11-14. 6
- [14] Sangyun Lee, Zinan Lin, and Giulia Fanti. Improving the training of rectified flows. arXiv preprint arXiv:2405.20320, 2024. 3
- [15] Baiqi Li, Zhiqiu Lin, Wenxuan Peng, Jean de Dieu Nyandwi, Daniel Jiang, Zixian Ma, Simran Khanuja, Ranjay Krishna, Graham Neubig, and Deva Ramanan. Naturalbench: Evaluating vision-language models on natural adversarial samples. arXiv preprint arXiv:2410.14669, 2024. 6
- [16] Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model. arXiv preprint arXiv:2311.06214, 2023. 2, 3, 4, 6, 7, 12
- [17] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 300–309, 2023. 2
- [18] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 2
- [19] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. Advances in Neural Information Processing Systems, 36, 2024. 3
- [20] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tok-makov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9298–9309, 2023. 3
- [21] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022. 2
- [22] Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, et al. Instaflow: One step is enough for high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference on Learning Representations*, 2023. 2, 3
- [23] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 2, 3

- [24] Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with pretrained models. Advances in Neural Information Processing Systems, 36, 2024.
- [25] Weijian Luo, Tianyang Hu, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and Zhihua Zhang. Diff-instruct: A universal approach for transferring knowledge from pre-trained diffusion models. In *NeurIPS*, 2023. 3
- [26] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In CVPR, 2023.
- [27] Paritosh Mittal, Yen-Chi Cheng, Maneesh Singh, and Shubham Tulsiani. AutoSDF: Shape priors for 3d completion, reconstruction and generation. In *CVPR*, 2022. 2
- [28] Thuan Hoang Nguyen and Anh Tran. Swiftbrush: One-step text-to-image diffusion model with variational score distillation. In CVPR, 2024. 3
- [29] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. arXiv preprint arXiv:2212.08751, 2022. 2, 3
- [30] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 6, 12
- [31] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv* preprint arXiv:2209.14988, 2022. 2, 6, 7
- [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 6
- [33] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021. 2
- [34] Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-guided texturing of 3d shapes. In *ACM SIGGRAPH 2023 conference proceedings*, pages 1–11, 2023. 3
- [35] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2
- [36] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. Advances in neural information processing systems, 35:36479–36494, 2022. 2
- [37] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 2, 3, 4

- [38] Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rombach. Fast highresolution image synthesis with latent adversarial diffusion distillation. arXiv preprint arXiv:2403.12015, 2024. 2, 3
- [39] Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen, Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base model. arXiv preprint arXiv:2310.15110, 2023.
- [40] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023. 2, 3, 4
- [41] Yawar Siddiqui, Antonio Alliegro, Alexey Artemov, Tatiana Tommasi, Daniele Sirigatti, Vladislav Rosov, Angela Dai, and Matthias Nießner. Meshgpt: Generating triangle meshes with decoder-only transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19615–19625, 2024. 3
- [42] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 4
- [43] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456, 2020. 2, 4
- [44] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *ICML*, 2023. 2, 3
- [45] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Viewset diffusion:(0-) image-conditioned 3d generative models from 2d data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8863– 8873, 2023. 2, 3
- [46] Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm: Large multi-view gaussian model for high-resolution 3d content creation. In *European Conference on Computer Vision*, pages 1–18. Springer, 2025. 6, 7, 12
- [47] Ayush Tewari, Tianwei Yin, George Cazenavette, Semon Rezchikov, Josh Tenenbaum, Frédo Durand, Bill Freeman, and Vincent Sitzmann. Diffusion with forward models: Solving stochastic inverse problems without direct supervision. Advances in Neural Information Processing Systems, 36: 12349–12362, 2023. 2, 3
- [48] Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, Adam Letts, Yangguang Li, Ding Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. Triposr: Fast 3d object reconstruction from a single image. *arXiv preprint arXiv:2403.02151*, 2024. 6, 12
- [49] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661– 1674, 2011.
- [50] Vikram Voleti, Chun-Han Yao, Mark Boss, Adam Letts, David Pankratz, Dmitrii Tochilkin, Christian Laforte, Robin Rombach, and Varun Jampani. SV3D: Novel multi-view synthesis and 3D generation from a single image using latent video diffusion. In European Conference on Computer Vision (ECCV), 2024. 3, 6, 12

- [51] Chen Wang, Jiatao Gu, Xiaoxiao Long, Yuan Liu, and Lingjie Liu. Geco: Generative image-to-3d within a second. arXiv preprint arXiv:2405.20327, 2024. 3
- [52] Peng Wang, Hao Tan, Sai Bi, Yinghao Xu, Fujun Luan, Kalyan Sunkavalli, Wenping Wang, Zexiang Xu, and Kai Zhang. Pf-lrm: Pose-free large reconstruction model for joint pose and shape prediction. arXiv preprint arXiv:2311.12024, 2023. 3
- [53] Tengfei Wang, Bo Zhang, Ting Zhang, Shuyang Gu, Jianmin Bao, Tadas Baltrusaitis, Jingjing Shen, Dong Chen, Fang Wen, Qifeng Chen, et al. Rodin: A generative model for sculpting 3d digital avatars using diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4563–4573, 2023. 2, 3
- [54] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. In *NeurIPS*, 2023. 2, 3
- [55] Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan, Valentin Deschaintre, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Meshlrm: Large reconstruction model for highquality mesh. arXiv preprint arXiv:2404.12385, 2024. 3
- [56] Shuang Wu, Youtian Lin, Feihu Zhang, Yifei Zeng, Jingxi Xu, Philip Torr, Xun Cao, and Yao Yao. Direct3d: Scalable image-to-3d generation via 3d latent diffusion transformer. arXiv preprint arXiv:2405.14832, 2024. 2, 3
- [57] Yinghao Xu, Hao Tan, Fujun Luan, Sai Bi, Peng Wang, Jiahao Li, Zifan Shi, Kalyan Sunkavalli, Gordon Wetzstein, Zexiang Xu, et al. Dmv3d: Denoising multi-view diffusion using 3d large reconstruction model. arXiv preprint arXiv:2311.09217, 2023. 2, 3
- [58] Yanwu Xu, Yang Zhao, Zhisheng Xiao, and Tingbo Hou. Ufogen: You forward once large scale text-to-image generation via diffusion gans. In CVPR, 2024. 3
- [59] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and William T Freeman. Improved distribution matching distillation for fast image synthesis. arXiv preprint arXiv:2405.14867, 2024. 2, 3, 4, 5
- [60] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6613–6623, 2024. 2, 3,
- [61] Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang Xu. Gs-lrm: Large reconstruction model for 3d gaussian splatting. In *European Conference on Computer Vision*, pages 1–19. Springer, 2025. 2, 3, 5, 8, 12
- [62] Longwen Zhang, Ziyu Wang, Qixuan Zhang, Qiwei Qiu, Anqi Pang, Haoran Jiang, Wei Yang, Lan Xu, and Jingyi Yu. Clay: A controllable large-scale generative model for creating high-quality 3d assets. ACM Transactions on Graphics (TOG), 43(4):1–20, 2024. 2, 3

## Turbo3D: Ultra-fast Text-to-3D Generation

# Supplementary Material

# 7. Details of Multi-step Multi-view Generation Model

We directly fine-tune an internal DiT [30] based text-to-image model into a text-to-multiview model. We fine-tune the model on the Objaverse dataset [3]. For the generation task, we render the dataset at a fixed elevation (20 degrees) and 16 equidistant azimuths. We empirically find that training with random views performs better than fixed views. In particular, during training, we randomly sample f views from the rendered 16 views for each instance, where f can be 4 or 8. Each view is conditioned on the corresponding Plücker embedding. For inference, we only infer 4 views for efficiency.

## 8. Experiments on 512 resolution

Some of the previous methods (Instant3D and SV3D) generate results with a higher resolution of 512. For a fair comparison, we also perform experiments on 512 resolution. Tab. 4 presents the quantitative comparisons with several state-of-the-art methods, where inference time is all measured under the resolution of 512. We can see our Turbo3D-512 version performs slightly better than our Turbo3D (256 resolution) with longer inference time, while outperforming other state-of-the-art methods by a large margin in terms of Clip score, VQA score, and inference speed.

Tab. 5 displays the effectiveness of latent GS-LRM. Under a higher resolution, the speed-up gain for latent GS-LRM gets larger. Overall, the latent GS-LRM archives a speed-up of 0.34s for the whole text-to-3D process.

## 9. Details of User Study

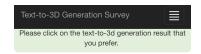
The interface example is shown in Fig. 7. For each question, we show two rendered videos from two different methods and ask the user to pick their preferred one. The two methods are randomly chosen from the total 4 methods: LGM [46], Instant3D [16], our multi-step multi-view model and our Turbo3D.

Method	CLIP	VQA	Inference
	Score ↑	Score ↑	Time ↓
TripoSR [48]	23.85	0.57	<b>1.28s</b> 35.96s
SV3D [50]	24.92	0.64	
Instant3D [16]	26.23	0.65	20.00s
LGM [46]	24.73	0.58	6.56s
Turbo3D-512	<b>27.66</b>	<b>0.78</b>	<b>1.28s</b>

Table 4. Comparison against state-of-the-art 3D generation methods. Our Turbo3D-512 generates 3D assets with highest CLIP and VQA scores while using the least amount of time (benchmarked on a A100 GPU).

Ablation	CLIP	VQA	Inference
	Score ↑	Score ↑	Time ↓
Pixel GS-LRM [61]	<b>27.68</b> 27.66	0.78	1.62s
Latent GS-LRM		0.78	<b>1.28s</b>

Table 5. Comparison between pixel and latent GS-LRM. We report the CLIP score, VQA score, and overall text-to-3D inference time for comparison. Our latent GS-LRM achieves similar image quality while enabling better efficiency (~21% speedup).



a tree stump with an axe buried in it



Method 1



Method 2



Figure 7. Interface example for user study.