

BOOTSTRAP3D: IMPROVING MULTI-VIEW DIFFUSION MODEL WITH SYNTHETIC DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent years have witnessed remarkable progress in multi-view diffusion models for 3D content creation. However, there remains a significant gap in image quality and prompt-following ability compared to 2D diffusion models. A critical bottleneck is the scarcity of high-quality 3D data with detailed captions. To address this challenge, we propose **Bootstrap3D**, a novel framework that automatically generates an arbitrary quantity of multi-view images to assist in training multi-view diffusion models. Specifically, we introduce a data generation pipeline that employs (1) 2D and video diffusion models to generate multi-view images based on constructed text prompts, and (2) our fine-tuned 3D-aware **MV-LLaVA** for filtering high-quality data and rewriting inaccurate captions. Leveraging this pipeline, we have generated 1 million high-quality synthetic multi-view images with dense descriptive captions to address the shortage of high-quality 3D data. Furthermore, we present a **Training Timestep Reschedule (TTR)** strategy that leverages the denoising process to learn multi-view consistency while maintaining the original 2D diffusion prior. Extensive experiments demonstrate that Bootstrap3D can generate high-quality multi-view images with superior aesthetic quality, image-text alignment, and maintained view consistency.

1 INTRODUCTION

3D content creation stands as a fundamental challenge within the generative domain, boasting widespread applications in augmented reality (AR) and game modeling. Unlike 2D image generation, the dearth of high-quality 3D models persists as a significant hurdle to overcome. In the realm of 2D image generation, the pivotal role of training on billion-scale image-text pairs (Schuhmann et al., 2022) has been firmly established (Betker et al., 2023; Rombach et al., 2022; Li et al., 2024; Chen et al., 2023a; 2024a). However, in 3D content generation, the scarcity of high-quality 3D models often compels reliance on the priors of 2D diffusion models. The predominant methodologies in this domain can be categorized into two main streams: 1) Gaining optimized neural representations from fixed 2D diffusion models via Score Distillation Sampling (SDS) loss (Poole et al., 2022; Shi et al., 2023b; Liu et al., 2023b; Shi et al., 2023a; Liu et al., 2023a; Wang et al., 2024a), which are time-intensive, lacking diversity and suffer from low robustness although capable of producing high-quality 3D objects. 2) Fine-tuning 2D diffusion models to achieve multi-view generation (Li et al., 2023a; Shi et al., 2023a;b), directly synthesizing 3D objects through sparse reconstruction models (Li et al., 2023a; Wang et al., 2023b; Xu et al., 2024a;b; Tang et al., 2024a; Wei et al., 2024). With recent improvements in large-scale sparse view reconstruction models and 3D representations (Kerbl et al., 2023), the second stream is garnering increasing attention.

Fine-tuning 2D diffusion models for multi-view generation remains challenging owing to the insufficiency in both data quality and quantity. Previous methods (Qiu et al., 2023; Li et al., 2023a; Shi et al., 2023b; Deitke et al., 2024) primarily train on a filtered subset of high-quality data from Objaverse (Deitke et al., 2023) and Objaverse-XL (Deitke et al., 2024). The scarcity of high-quality data often introduces various shortcomings. In single-view based novel view synthesis (Liu et al., 2023b; Shi et al., 2023a; Wang & Shi, 2023; Voleti et al., 2024), if the input images deviate from the distribution of the training data, it can induce issues such as motion blurring, object distortion and deformation (Shi et al., 2023a).



085
086
087

Figure 1: **Bootstrap3D** can generate high quality multi-view images with precise long text control and style customization while maintaining view consistency.

088
089
090
091
092
093
094

Moreover, in direct text-to-multi-view image generation, the pursuit of enhancing view consistency compromises the aesthetic and photo-realistic quality. For instance, Instant3D (Li et al., 2023a) fine-tunes SDXL (Podell et al., 2023) using only 10K high-quality Objaverse (Deitke et al., 2023) data with a small learning rate for 10K steps, which does not fundamentally prevent the catastrophic forgetting problem of losing 2D diffusion prior, leading to compromised image quality. Recent endeavors have predominantly focused on alleviating data scarcity and improving view consistency from a model-centric perspective (Kant et al., 2024; Shi et al., 2023a; Tang et al., 2024b), with limited exploration into the improvement of training data and training method itself.

095
096
097
098
099
100
101
102
103
104
105
106
107

Recent Multimodal Large Language Models (MLLMs) (Liu et al., 2024a; Chen et al., 2023b; Li et al., 2023b; Alayrac et al., 2022; Anil et al., 2023) like GPT-4V (OpenAI, 2023a) and Gemini (Team et al., 2023), possess image understanding capabilities and rudimentary 3D world awareness, has enabled automatic quality assessment of multi-view images and dense caption generation. Furthermore, notable advancements in video diffusion (Brooks et al., 2024; Voleti et al., 2024) have improved the generalizability of novel view synthesis (Voleti et al., 2024; Chen et al., 2024b; Kwak et al., 2023). Employing these advancements, we propose Bootstrap3D to generate synthetic data to counteract the data deficiencies inherent in training multi-view diffusion models. To be specific, we introduce the Bootstrap3D data generation pipeline for producing high-quality multi-view images with dense descriptive captions. Subsequently, we fine-tune a multi-view-aware MLLM model, dubbed as MV-LLaVA, to achieve fully automated high-quality data annotation with both efficiency and accuracy. To mitigate catastrophic forgetting of 2D diffusion prior, we introduce a training timestep reschedule (TTR) strategy when fine-tuning multi-view diffusion models. Specifically, we use the phased nature of the denoising process (Ho et al., 2020) and limit different training time steps for synthetic data to achieve enhanced image quality with maintained view consistency.

108 Through extensive experiments, we demonstrate that our method significantly enhances the adher-
109 ence of the multi-view diffusion model to text prompts and image quality while ensuring view con-
110 sistency. Integrated with the reconstruction model, our approach facilitates the creation of 3D models
111 with superior quality. We show some of the qualitative results in Fig. 1, where our model can achieve
112 high quality multi-view images with precise text control and style customization. Our contributions
113 are summarized into the following points:

114 **1)** We present an automated **Bootstrap3D** data generation pipeline that uses the video diffusion
115 model and our fine-tuned 3D-aware MV-LLaVA to synthesize an arbitrary number of high-quality
116 multi-view image text pairs.

117 **2)** We propose a Training Time-step Reschedule (TTR) strategy for fine-tuning the multi-view diffu-
118 sion model that employs both synthetic data and real data to enhance image quality and image-text
119 alignment while maintaining view consistency.

120 **3)** We generate 1 million multi-view images with dense descriptive captions suitable for training
121 the multi-view diffusion model and provide dense descriptive captions on Objaverse Deitke et al.
122 (2023), which mitigates the gap with the 2D diffusion model from a data perspective.

125 2 RELATED WORK

127 **Existing 3D datasets and data pre-processing.** Existing object level 3D datasets, sourced either
128 from CAD (Chang et al., 2015; Wu et al., 2015; Deitke et al., 2023; 2024) or scan from real ob-
129 jects (Aanæs et al., 2016; Yao et al., 2020; Downs et al., 2022; Wu et al., 2023), are still small in size.
130 Most state-of-the-art open-sourced 3D content creation models are trained on Objaverse (Deitke
131 et al., 2023). However, there still exists a huge gap compared to data used for training 2D diffusion
132 models (Schuhmann et al., 2022). In addition to quantity, quality is also an important problem re-
133 mains to be solved as many methods (Shi et al., 2023b; Li et al., 2023a; Qiu et al., 2023; Tang et al.,
134 2024a) trained on Objaverse rely on filtering out low-quality data, making the precious 3D data even
135 less. Another critical gap that requires attention is the quality of the 3D object’s caption. Previous
136 work Cap3D (Luo et al., 2024) propose to apply BLIP-2 (Li et al., 2023b) and GPT-4 (OpenAI,
137 2023b) to generate caption based on multi-view images. However, this approach, without direct
138 input image into GPT, can lead to severe hallucination. Given recent breakthroughs in improving
139 text-image alignment through caption rewriting (Betker et al., 2023; Chen et al., 2023a; 2024a; Esser
140 et al., 2024), there is a pressing need to rewrite denser and more accurate captions for 3D objects
141 with the assistance of advanced Multimodal Large Language Models (MLLMs). In this work, we
142 propose a new data generation pipeline to synthesize multi-view images and rewrite captions for 3D
143 objects incorporating additional quality scoring mechanisms to address the aforementioned issues.

144 **Text-to-3D content creation.** The field of 3D content creation has been a vibrant area of research
145 over the past years. One prominent research direction explores the use of Score Distillation Sam-
146 pling (SDS) (Poole et al., 2022) and its variants (Chen et al., 2023c; Chung et al., 2023; Hertz et al.,
147 2023; Liang et al., 2023; Lin et al., 2023; Liu et al., 2023b; Shi et al., 2023b; Liu et al., 2023c; Long
148 et al., 2023; Wang et al., 2024a; Tang et al., 2023; Wang et al., 2023a; Yang et al., 2024; Qi et al.,
149 2024), using the priors of 2D diffusion models to optimize 3D representations. While these methods
150 have demonstrated success in producing high-quality 3D generations, they often require prolonged
151 optimization time to converge. In contrast, recent studies (Hong et al., 2023; Wang et al., 2023b;
152 Li et al., 2023a; Tang et al., 2024a; Tochilkin et al., 2024; Xu et al., 2024b; Wei et al., 2024) have
153 proposed the direct inference of 3D representations (Mildenhall et al., 2021; Chan et al., 2022; Kerbl
154 et al., 2023; Zhang et al., 2023a) conditioned by images. Among these approaches, Instant3D (Li
155 et al., 2023a) stands out by utilizing multi-view images of the same object to directly deduce the Tri-
156 plane (Chan et al., 2022) representation. This approach effectively addresses the issue of ambiguous
157 unseen areas inherent in the single image to 3D conversions, as encountered in LRM (Hong et al.,
158 2023) and TripoSR (Tochilkin et al., 2024). Instant3D, along with subsequent works (Xu et al.,
159 2024b; Zheng et al., 2024; Wang et al., 2024b; Xu et al., 2024a), efficiently decomposes 3D genera-
160 tion into two processes: the generation of multi-view images using multi-view diffusion model (Liu
161 et al., 2023b;c;a; Shi et al., 2023b; Liu et al., 2024b; Shi et al., 2023a; Long et al., 2023; Kant et al.,
2024; Voleti et al., 2024) and large reconstruction model to generate 3D representations conditioned
on these multi-view images. In this work, we introduce a method that significantly enhances the
scalability of training and data generation for multi-view image generation.

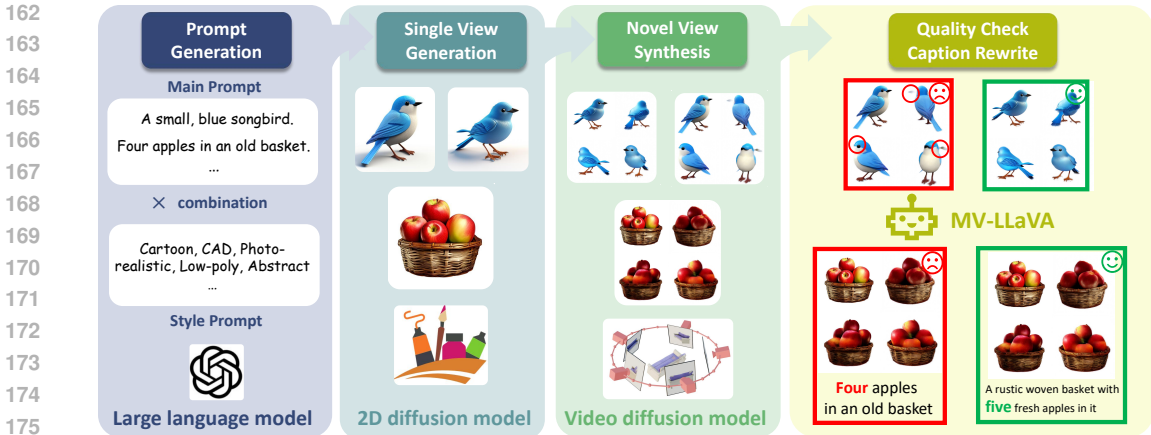


Figure 2: **Bootstrap3D data generation pipeline** that consists of 1) using LLM to generate diverse text prompts 2) employing the T2I model to generate single-view images 3) synthesizing arbitrary number of multi-view images by applying the video diffusion model, 4) employing MV-LLaVA to filter and select only high-quality data, and rewrite captions to be dense and descriptive.

Video diffusion for novel view synthesis. Recent advancements in video diffusion have marked a significant breakthrough, with models such as Sora (Brooks et al., 2024) and SVD (Blattmann et al., 2023) scaling up the direct generation process from images to videos. Following these developments, a series of works (Wang et al., 2023c; Kwak et al., 2023; Blattmann et al., 2023; Melas-Kyriazi et al., 2024; Han et al., 2024; Chen et al., 2024b) represented by SV3D (Voleti et al., 2024), have fine-tuned these video diffusion models for 3D content creation. Despite these groundbreaking developments, the new perspective images generated based on video priors still suffer from issues like motion blur. In this work, we propose to utilize SV3D (Voleti et al., 2024) as a data generator to produce novel views of given images with additional quality checks to leave only high-quality data.

Multimodal Large Language Models. With the development of large language models (Brown et al., 2020; OpenAI, 2023b; Chowdhery et al., 2022; Anil et al., 2023; Hoffmann et al., 2022; Touvron et al., 2023), multimodal large language models (MLLMs) (Zhang et al., 2023b; Alayrac et al., 2022; Li et al., 2023b; 2022; Huang et al., 2023; Driess et al., 2023; Awadalla et al., 2023; Liu et al., 2024a; Dong et al., 2024; Sun et al., 2023), such as GPT-4V (OpenAI, 2023a), have demonstrated groundbreaking 2D comprehension capabilities and open-world knowledge. As is discovered in GPTEval3D (Wu et al., 2024), GPT-4V can achieve human-aligned evaluation for multi-view images rendered from 3D objects. In this work, we fine-tune the LLaVA (Liu et al., 2024a) for quality judgment and descriptive caption generation based on multi-view images.

3 METHODS

Due to the scarcity of high-quality 3D data, we develop the Bootstrap3D data generation pipeline to efficiently construct an arbitrary number of training data (Sec. 3.1). Subsequently, the quality of generated multi-view images is assessed using the powerful GPT-4V (OpenAI, 2023a) or our proposed MV-LLaVA (Liu et al., 2024a) model to generate dense descriptive captions efficiently and faithfully (Sec. 3.2). We also design a training timestep reschedule (Sec. 3.3) when fine-tuning the multi-view diffusion model with our synthetic and real data.

3.1 BOOTSTRAP3D DATA GENERATION PIPELINE

As illustrated in Fig.2, our data generation pipeline initially employs GPT-4 (OpenAI, 2023a) to generate a multitude of imaginative and varied text prompts (Wu et al., 2024). Subsequently, to generate 2D images that closely align with the text prompts, we utilize the PixArt-Alpha (Chen et al., 2023a) model use FlanT5 (Chung et al., 2024) text encoder with DiT (Peebles & Xie, 2023) architecture for text-to-image (T2I) generation. Thereafter, we use SV3D (Voleti et al., 2024) for novel view synthesis. Given the significant motion blur and distortion often present in SV3D (Voleti et al., 2024) outputs, we further employ Multimodal Large Language Models(MLLM) to evaluate the quality of multi-view images. To rectify mismatches between multi-view images and the original

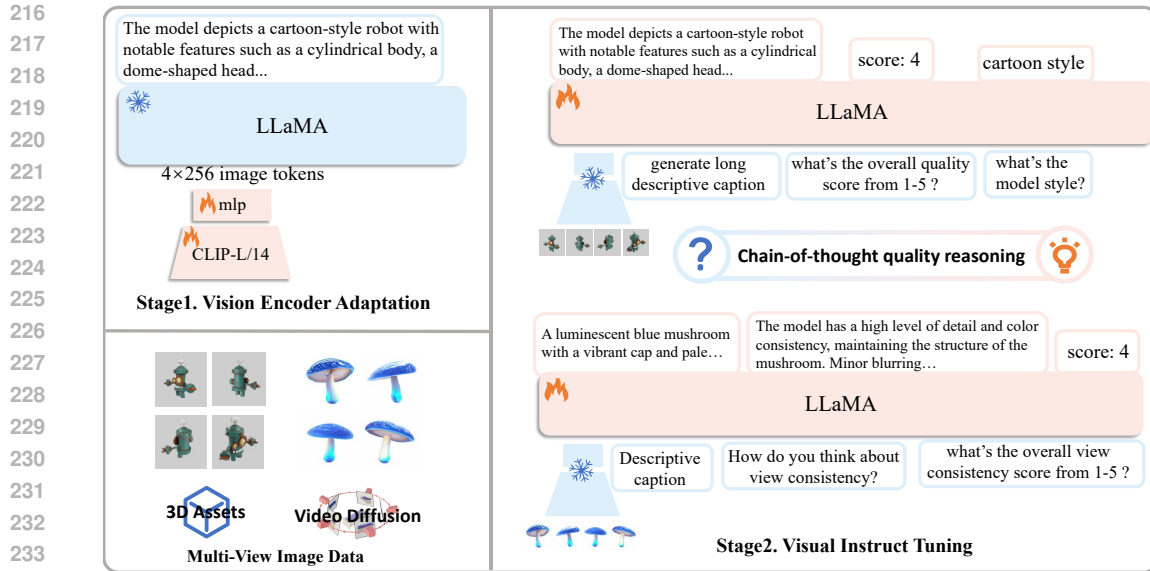


Figure 3: **MV-LLaVA**. We use GPT-4V (OpenAI, 2023a) to generate long descriptive captions, quality scoring, and reasoning processes for multi-view images to construct the instruction tuning dataset. Then we fine-tune our MV-LLaVA based on LLaVA (Liu et al., 2024a) to serve as the human-aligned quality checker and captioner for multi-view images.

text prompts induced by novel view synthesis and provide more precise captions, we further propose MV-LLaVA to generate dense descriptive captions for multi-view images.

3.2 MULTI-VIEW LLaVA (MV-LLaVA)

To efficiently generate captions and label quality scores for both generated multi-view images and 3D assets in Objaverse (Deitke et al., 2023), we propose the Multi-View LLaVA (MV-LLaVA) that fine-tune LLaVA (Liu et al., 2024a) based on our instructive conversation pairs generated by the powerful GPT-4V (OpenAI, 2023a).

Preparing the instruction tuning data. As shown in Fig.2, we use GPT-4 to generate 20k varied text prompts based on prompts designed in (Wu et al., 2024) and use PixArt-alpha (Chen et al., 2023a) to generate single view image and use SV3D (Voleti et al., 2024) or Zero123++ (Shi et al., 2023a) to generate multi-view images. For these 20k generated multi-view images, we prompt GPT-4V (OpenAI, 2023a) to generate comments on view consistency, image quality and generate dense descriptive captions. For the additional 10K rendered multi-view images from Objaverse (Deitke et al., 2023), we prompt GPT-4V (detailed prompts in Sup. A.5.1) to offer feedback on the quality and aesthetic appeal of 3D objects, along with style judgments. We utilize these 30K high-quality multi-view image text pairs (prompts detailed in Sup. A.5.2) as the instruction tuning data for LLaVA.

Instruction tuning. As presented in the left part of Fig. 3, due to the LLaVA’s maximum training context length constraints of 2048, we input four images separately into CLIP-L/14 (Radford et al., 2021) and generate 4×256 image tokens. Inspired by ShareGPT-4V (Chen et al., 2023b), we freeze only a portion of layers of CLIP (Radford et al., 2021) in the first stage of pre-training to enhance multi-view awareness and texture perception of vision encoder (detailed in Sup. A.4.1). As shown in the right part of Fig. 3, we first ask the model to generate descriptions, then let the model score the quality based on multi-view images and captions. Our approach encourages LLM to deduct more reasonable quality scores through chain-of-thought (Wei et al., 2022). A mixture of original training data of LLaVA is adopted to mitigate over-fitting. As a result, we obtain MV-LLaVA, which efficiently filters and re-captions both synthetic data and 3D assets. As detailed in Sup.A.4, MV-LLaVA can not only generate more accurate, less hallucinated dense captions that faithfully describe 3D objects compared to Cap3D (Luo et al., 2024) but also assign the human-aligned quality score on both synthetic data and Objaverse assets. The filtered high-quality multi-view images with rewritten dense captions served as training data for the diffusion model.

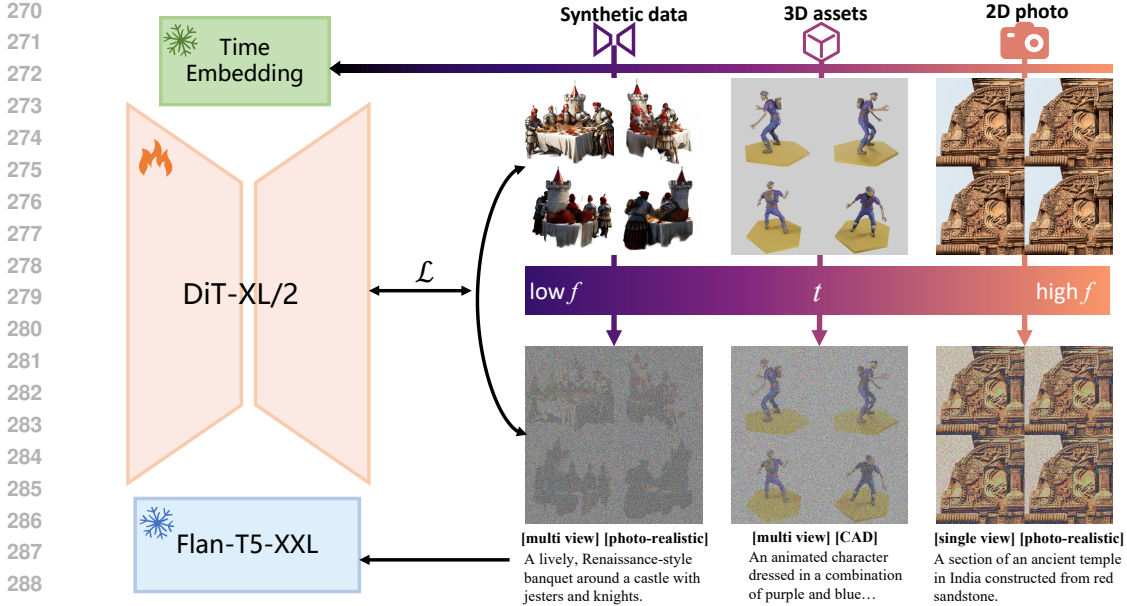


Figure 4: **Training Timestep Reschedule (TTR)**. For different types of training data, we restrict the training time step t accordingly to achieve the balance between varied high aesthetic images that are better aligned with text prompt, photo-realistic texture, and view consistency for 3D generation.

3.3 TRAINING TIMESTEP RESCHEDULE (TTR)

Despite retaining only relatively high-quality synthetic data with minimal motion blur from SV3D (Voleti et al., 2024) through MV-LLaVA, small areas of blurring persist, stemming from both motion and out-of-distribution scenarios for SV3D and SVD (Blattmann et al., 2023). These blurred data can potentially compromise the final performance of the multi-view diffusion model. To restrict the training time step for synthetic data, we proposed a simple yet effective Training Timestep Reschedule (TTR) method.

Background. Before delving into TTR, we briefly review some basic concepts needed to understand diffusion models (DDPMs) (Ho et al., 2020; Sohl-Dickstein et al., 2015; Salimans & Ho, 2022; Rombach et al., 2022; Chen et al., 2023a). Gaussian diffusion models assume a forward noising process which gradually applies noise to real data x_0

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (1)$$

here constants $\bar{\alpha}_t$ are hyperparameters. By applying the reparameterization trick, we can sample

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

During training, t is randomly sampled in $[0, N]$ ($N = 1000$ in (Chen et al., 2023a; Rombach et al., 2022)) for the model to predict the added noise ϵ_t , where x_0 denotes for the clear nature image and x_N denotes for pure Gaussian noise. As depicted in Fig.4, when t is large, the denoising process primarily focuses on determining the global low frequency(f) content such as overall structure and shape. Conversely, when t is small, the denoising process is predominantly responsible for generating high f components such as texture.

When adapting Stable Diffusion (Rombach et al., 2022) for multi-view generation, the previous approach (Shi et al., 2023a) changes the default scaled linear schedule into the linear schedule to emphasize more on early denoising stage for structural variation and view consistency. Inspired by this, we propose restricting the denoising time step of synthetic data during training. As small yet observable blur still exists in synthetic data with novel view generated by SV3D (Voleti et al., 2024), we limit them to training diffusion model only with large t . This restricts the backpropagation of these synthetic data to focus on the low-frequency component of the image like the overall structure and shape that faithfully follow text prompts and consistency between different views. Small t values are only sampled on clear and physically consistent multi-view images rendered from Obja-verse (Deitke et al., 2023) and supplemented high-quality 2D images from SA-1B (Kirillov et al.,



351 Figure 5: **Bootstrap3D generates 3D objects compared to other edge-cutting methods** given text
352 prompt. More results with higher resolution are available in Sup.A.8.1.

353 2023), help model outcome high-quality images with more photo-realistic and varied texture details.
354

356 4 EXPERIMENTS

358 4.1 EXPERIMENT SETTINGS

360 **Training data.** For each set of 4-view images obtained from both Objaverse (Deitke et al., 2023)
361 and generated by SV3D (Voleti et al., 2024) or Zero123++(Shi et al., 2023a), we use MV-LLaVA
362 to generate long descriptive captions with predicted quality score. Detailed quality check of MV-
363 LLaVA is supplied in Sup. A.4 and data analysis in Sup. A.3. In the end, we generate 200K 4-view
364 image-text pairs on Objaverse (Deitke et al., 2023), 1000K 4-view image-text pairs from synthetic
365 data from SV3D (Voleti et al., 2024) and Zero123++(Shi et al., 2023a). We also sample 35K HQ
366 SA (Kirillov et al., 2023) data with captions from ShareGPT4V (Chen et al., 2023b).

367 **Training details.** We test our framework directly on the text-to-multi-view diffusion model. We
368 fine-tune PixArt- α (Chen et al., 2023a) with backbone DiT-XL/2 (Peebles & Xie, 2023) model
369 on the data as mentioned earlier. Similar to Instant3D (Li et al., 2023a), we train the diffusion
370 model directly on 4-view images naturally arranged in a 2×2 grid. For 4 same view images from
371 SA (Kirillov et al., 2023), we limit training time step $t \in [0, 50]$. We limit synthetic multi-view
372 images $t \in [200, 1000]$. Regarding 3D object-rendered images, we do not limit t but sample more
373 frequently in the range $[50, 200]$ as a complement. We set the total batch size to 1024 with the
374 learning rate set to $8e-5$ for 20K steps. Training is conducted on 32 NVIDIA A100-80G GPUs for
375 20 hours with Flan-T5-XXL (Chung et al., 2024) text features and VAE (Kingma & Welling, 2013)
376 features pre-extracted.

377 **Evaluation metrics.** We primarily benchmark the quantitative results of our approach and other
methods from two main dimensions: 1). **Image-text alignment** measured by CLIP score and CLIP-

Table 1: **Benchmark of CLIP and FID score of text-to-multi-view (T2MV) models** on generated 4 view images, CLIP score tests on 110 text prompts from GPTeval3D (Wu et al., 2024) while FID is measured with the distribution of 30K object-centric images generated by SOTA T2I models. For text-to-image-to-multi-view(T2I2MV), we input I2MV models with single view images generated by Pixart- α , which superior single view image quality is marked in green .

Domain	Method	CLIP-R Score \uparrow		CLIP Score \uparrow		FID \downarrow	
		CLIP-L/14	CLIP-bigG	CLIP-L/14	CLIP-bigG	PG2.5	PixArt- α
T2I	PixArt- α	96.1	94.7	25.9	41.5	20.7	5.4
T2I2MV	SV3D	78.8	81.3	24.7	37.3	55.7	54.2
	CRM	77.5	85.1	24.9	38.9	59.0	52.2
	Zero123++	78.0	84.5	24.2	36.9	53.2	49.3
T2MV	Instant3D (unofficial)	83.6	91.1	25.6	39.2	83.2	77.9
	MVDream	84.8	89.3	25.5	38.4	60.2	59.2
	Bootstrap3D	88.8	92.5	25.8	40.1	42.4	31.0

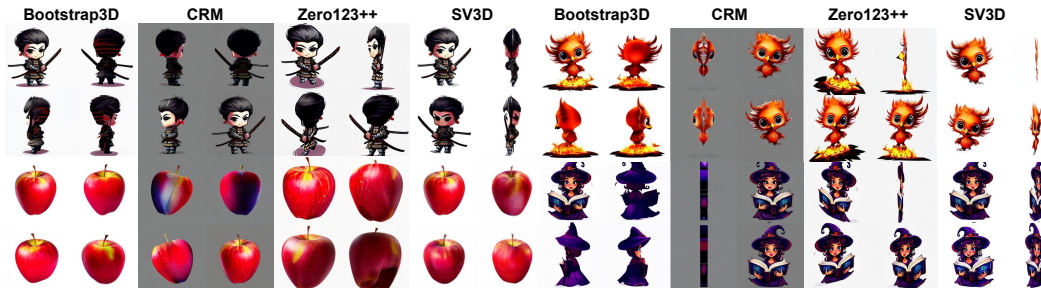


Figure 6: **Bootstrap3D can generate high quality multi-view images** in out of domain cases compare to other edge-cutting multi-view diffusion models trained on Objaverse only.

R score indicating the prompt follow ability of text-to-multi-view (T2MV) diffusion model. 2). **Quality of generated images** measured by FID (Heusel et al., 2017). Given the trend of decoupling multi-view image generation and sparse view reconstruction, we conduct tests separately on multi-view images by T2MV and rerendered images from generated 3D objects. To test the robustness and diversity of Bootstrap3D beyond prompts generated by GPT, we also collect real user prompts from public website, the details and test results are available in Sup. A.1.

Evaluation details. For CLIP-R Score and CLIP Score, we test on 110 text prompts from GPTeval3D (Wu et al., 2024) using different CLIP models (i.e., CLIP-L/14 (Radford et al., 2021) and CLIP-bigG (Ilharco et al., 2021)) following the same setting of Instant3D (Li et al., 2023a). Regarding the FID (Heusel et al., 2017) test, as there is no golden standard for HQ 3D objects, we follow the similar evaluation idea of Playground2.5 (Li et al., 2024) (PG2.5) to use powerful T2I model generated images to form ground truth (GT) distribution. We use curated prompts to guide powerful PixArt and PG2.5 to generate high-quality CAD-style images with a single object in the pure background. Rembg (etc, 2020) is adopted to create white background object-centric images. We use the method proposed in GPTeval3D (Wu et al., 2024) to generate 3K prompts. For both PG-2.5 and PixArt, we generate 10 images for each prompt with different seeds, resulting in 30K images to form the GT distribution of high-quality CAD-style objects.

Comparing methods. In addition to Instant3D (Li et al., 2023a) and MVDream (Shi et al., 2023b) as direct text-to-multi-view (T2MV) methods, we also adopt edge-cutting single image to multi-view (I2MV) methods CRM (Wang et al., 2024b), SV3D (Voleti et al., 2024) and Zero123++(Shi et al., 2023a). For these methods, we condition the diffusion model on the single view image generated by PixArt (prompted to generate CAD-style single object-centric image). The result of the CLIP score is 3 times averaged with different seeds. For FID, we use 3 different seeds for each of the 3K prompts to generate 9K images to test the distance with GT high-quality images.

4.2 EVALUATION OF MULTI-VIEW IMAGES

As illustrated in Tab.1, compared to other methods, the T2MV diffusion model trained by our framework yields the best results both according to image-text alignment and image quality. For qualitative experiments, we visualize some of the comparisons with other edge-cutting multi-view diffusion

model in Fig.6. For these image-to-multi-view models, we condition them on the top-left image generated by Bootstrap3D. Compared to these models trained solely on Objaverse Deitke et al. (2023), our model demonstrates superior generalizability when the image domain is beyond the domain of Objaverse. Since it is difficult to directly measure view consistency as there is no ground truth 3D object for text-to-3D generation, we evaluate the view consistency by synthesizing 3D objects through large reconstruction model in the following experiments. Qualitative results of real user cases are in Sup. A.1.

4.3 EVALUATION OF GENERATED 3D OBJECTS

Table 2: **Benchmark of CLIP and FID score of generated 3D objects** based on rendered 9 view images. *MVDream is tested on 200 generated objects for FID test using SDS (Shi et al., 2023b), other methods are tested on 1000 objects using GRM (Xu et al., 2024b) and InstantMesh (Xu et al., 2024a) as sparse view reconstruction model.

Reconstruction	Method	CLIP-R Score \uparrow		CLIP Score \uparrow		FID \downarrow	
		CLIP-L/14	CLIP-bigG	CLIP-L/14	CLIP-bigG	PG2.5	PixArt
SDS	MVDream*	85.2	90.8	26.1	39.4	57.8	56.7
	Instant3D (unofficial)	81.7	89.4	24.8	37.1	85.4	80.3
GRM	SV3D	74.1	82.8	23.4	34.1	68.4	69.1
	Zero123++	71.2	80.3	22.3	34.5	69.3	72.4
	Bootstrap3D	86.3	91.6	25.9	39.7	51.2	50.7
	Zero123++	73.2	84.1	23.0	37.2	82.3	88.8
InstantMesh	Zero123++	73.2	84.1	23.0	37.2	82.3	88.8
	Bootstrap3D	87.1	92.0	26.0	39.2	61.2	55.3

View consistency is another crucial factor in reconstructing reasonable 3D objects. Miss alignment between different views can lead to blurred areas in reconstructed objects by large reconstruction model (Hong et al., 2023; Wei et al., 2024). This misalignment causes a significant deterioration in quality, resulting in a notable increase in metrics like FID. To assess the view consistency directly on 3D object, we employ GRM (Xu et al., 2024b) and InstantMesh (Xu et al., 2024a) to reconstruct the object given sparse view images generated in Sec. 4.2. We render 9 view images evenly in orbit for each object and evaluate the image-text alignment and image quality. As reported in Tab. 2, Bootstrap3D, after conditioning GRM or InstantMesh on 4 view images, can generate the best 3D objects both according to image-text alignment and image quality. GPT-4V based human-aligned evaluation based on GPTeval3D (Wu et al., 2024) is supplied in Sup. A.6.

We also present visualizations of some results in Fig.5. Bootstrap3D can generate objects with higher quality and prompt following ability. For other methods, as shown in the first column of Fig.5, although the first image may be well aligned with the given text prompt, the final 3D object may be compromised due to the limitations of its poor generalizability as they are also fine-tuned on Objaverse (Deitke et al., 2023) only. More visualizations and discussions of this are in Sup. A.2

4.4 ABLATION STUDY

Training Timestep Reschedule (TTR) is proposed in 3.3 to better integrate different types of data. The training time step of synthetic data is restricted in $[T, 1000]$, where T is a hyper-parameter to be set in training. We demonstrate the effect of the time-step limit in Fig.7, where the bar in the middle is the value of T . When T is large, namely synthetic data won't affect more time-step at the end of the denoising process, Synthetic data has less influence on the denoising process towards the end, which leads to better view consistency but lower prompt-following ability. Conversely, if T is small, the denoised result better follows the given text prompt but blurring becomes much more severe. In summary, there is a trade-off in injecting synthetic data into the training process: better image-text alignment comes at the cost of worse view consistency and increased blurring. Ultimately, we set $T = 200$ based on empirical study.

Synthetic data and dense captioning are proposed in our work to achieve high-quality images and better image-text alignment. We ablate their effects and the importance of data quantity in Tab. 3. Direct use of synthetic data without Training Timestep Reschedule (TTR) can cause severe blurs and deformation in final outcome. With the help of TTR, the mixture of data can not only improve image-text alignment but also maintain view consistency. Replacing Cap3D (Luo et al., 2024)'s caption

Table 3: **Ablation study of proposed components and quantity of synthetic data.** with CLIP-R Score represents image-text alignment and FID represents image quality.

Methods	Multi-view Image		Generated Object	
	CLIP-R Score	FID PG-2.5	CLIP-R Score	FID PG-2.5
Instant3D (unofficial)	83.6	83.2	81.7	85.4
Cap3D only	77.9	101.3	74.6	120.4
Cap3D + Synthetic Image (100k) w/o TTR	81.5	92.0	71.2	134.6
Cap3D + Synthetic Image (100k) w/ TTR	83.3	60.8	80.2	70.6
Dense recaption + Synthetic Image (100k)	87.4	50.2	85.1	50.9
Dense recaption + Synthetic Image (500k)	88.8	42.4	86.3	51.2

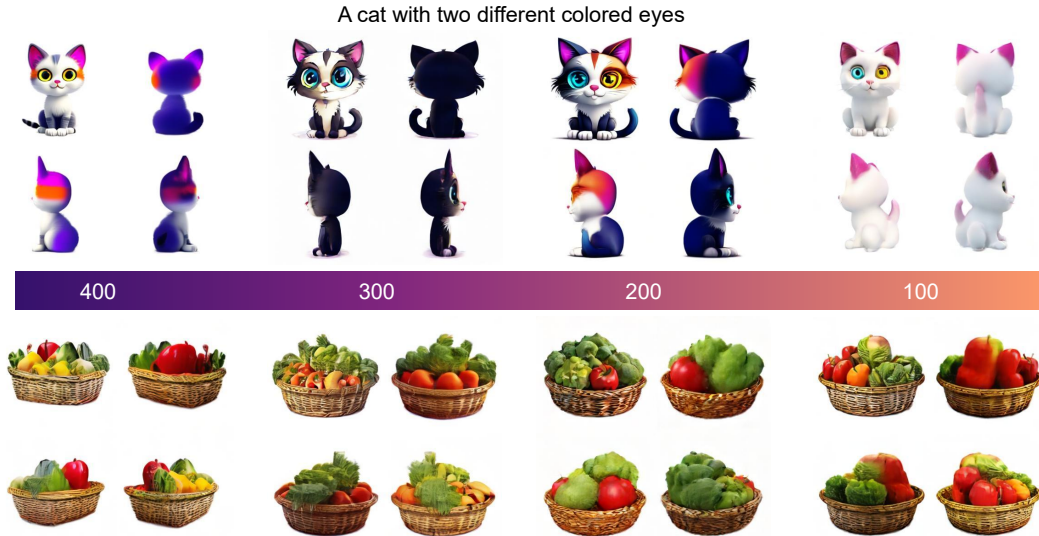


Figure 7: **Ablation study of training time reschedule (TTR)** demonstrates a trade-off between image-text alignment and image quality with different t .

with MV-LLaVA’s dense descriptive caption further improves the model’s capability of following prompts faithfully. Improvement through increasing volume of data also proves the scalability of our framework.

5 CONCLUSION AND DISCUSSION

In this work, we introduce a novel framework that employs MLLMs and diffusion models to synthesize high-quality data for bootstrapping multi-view diffusion models. With a powerful fine-tuned 3D-aware MLLM serving as the dense captioner and quality filter, the generated synthetic data addresses the issue of insufficient high-quality 3D data. The proposed strategy of injecting different data at different training time steps uses the property of the denoising process to further achieve higher image quality while maintaining view consistency. We believe this work will contribute to the goal of achieving 3D content creation with each rendered view comparable with the single view diffusion model, with more advanced MLLMs and diffusion models on the horizon.

Limitations and future work. Despite its promise, our work still faces several unresolved challenges. Firstly, the multi-view diffusion model is only the first step of the 3D content creation pipeline. Sparse view reconstruction models also need improvement as most edge-cutting sparse view reconstruction models are also trained on Objaverse Deitke et al. (2023) only. Secondly, Although MLLMs can estimate general quality and view consistency, subtle view inconsistency is hard to detect until ambiguity leads to blurred areas in reconstructed 3D object. While the proposed Training Timestep Reschedule can mitigate this problem, it cannot solve the problem fundamentally. Using synthetic data to train sparse view reconstruction models and quality estimation directly based on the reconstructed object are thus interesting future directions for improving 3D content creation.

6 ETHICS STATEMENT

Our training code is modified based on public available repository <https://github.com/PixArt-alpha/PixArt-alpha>. Part of training data are synthesized by our proposed data generation pipeline. For other part of original Objaverse Deitke et al. (2023) data, we only use Cap3D Luo et al. (2024) filtered assets (Objects with CC BY-NC-SA and CC BY-NC licenses are removed, while we retain those with CC-BY 4.0, CC BY-SA, and CC0 licenses) and with face recognizable objects filtered through MSFW classifier and face detector. The ethical filtering in Cap3D make our work using only data without ethics problem. For our synthetic new data, We will launch both the generated captions for Objaverse Deitke et al. (2023) and high-quality synthetic data, model checkpoints and codes with CC-BY 4.0 license for the research community.

7 REPRODUCIBILITY STATEMENT

Main experimental setting/details (training data, hyperparameters, optimizer, evaluation settings, etc) are clearly presents in Sec. 4.1. For main results, we detail the full test settings in Sec. 4.1. For GPT-4V OpenAI (2023a) based preference study, we provide detailed test prompts and test settings in Sup. A.6. Readers can easily follow the same settings and reproduce all of our experiment results. We provide code for generating synthetic data. Both codes for training the model and testing are also available in supplementary material. The full data and model checkpoints are too large to provide public link without violation of double-blinding. We will release full data and model checkpoints after review.

REFERENCES

- Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision*, 120: 153–168, 2016.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. *ArXiv*, abs/2204.14198, 2022. URL <https://api.semanticscholar.org/CorpusID:248476411>.
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Tachard Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Z. Chen, Eric Chu, J. Clark, Laurent El Shafey, Yanping Huang, Kathleen S. Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omer-nick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan A. Botha, James Bradbury, Siddhartha Brahma, Kevin Michael Brooks, Michele Catasta, Yongzhou Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, C Crépy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, M. C. D’iaz, Nan Du, Ethan Dyer, Vladimir Feinberg, Fan Feng, Vlad Fienber, Markus Freitag, Xavier García, Sebastian Gehrmann, Lucas González, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, An Ren Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wen Hao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Mu-Li Li, Wei Li, Yaguang Li, Jun Yu Li, Hyeontaek Lim, Han Lin, Zhong-Zhong Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pel-lat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alexandra Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Marie Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniela Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wiet-ing, Yuhuai Wu, Ke Xu, Yunhan Xu, Lin Wu Xue, Pengcheng Yin, Jiahui Yu, Qiaoling Zhang, Steven Zheng, Ce Zheng, Wei Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. Palm 2 techni-

- 594 cal report. *ArXiv*, abs/2305.10403, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:258740735)
595 [CorpusID:258740735](https://api.semanticscholar.org/CorpusID:258740735).
596
- 597 Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani
598 Marathe, Yonatan Bitton, Samir Yitzhak Gadre, Shiori Sagawa, Jenia Jitsev, Simon Korn-
599 blith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Open-
600 flamingo: An open-source framework for training large autoregressive vision-language mod-
601 els. *ArXiv*, abs/2308.01390, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:261043320)
602 [CorpusID:261043320](https://api.semanticscholar.org/CorpusID:261043320).
- 603 James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang
604 Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer*
605 *Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3), 2023.
606
- 607 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik
608 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling
609 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- 610 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe
611 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video
612 generation models as world simulators. 2024. URL [https://openai.com/research/](https://openai.com/research/video-generation-models-as-world-simulators)
613 [video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 614 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
615 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-
616 wal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh,
617 Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler,
618 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-
619 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot
620 learners. *ArXiv*, abs/2005.14165, 2020. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:218971783)
621 [CorpusID:218971783](https://api.semanticscholar.org/CorpusID:218971783).
- 622 Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio
623 Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d
624 generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision*
625 *and pattern recognition*, pp. 16123–16133, 2022.
- 626 Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li,
627 Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d
628 model repository. *arXiv preprint arXiv:1512.03012*, 2015.
629
- 630 Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James
631 Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for photore-
632 alistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023a.
- 633 Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping
634 Luo, Huchuan Lu, and Zhenguo Li. Pixart- σ : Weak-to-strong training of diffusion transformer
635 for 4k text-to-image generation. *arXiv preprint arXiv:2403.04692*, 2024a.
636
- 637 Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua
638 Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint*
639 *arXiv:2311.12793*, 2023b.
- 640 Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and
641 appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF Inter-*
642 *national Conference on Computer Vision*, pp. 22246–22256, 2023c.
- 643 Zilong Chen, Yikai Wang, Feng Wang, Zhengyi Wang, and Huaping Liu. V3d: Video diffusion
644 models are effective 3d generators. *arXiv preprint arXiv:2403.06738*, 2024b.
645
- 646 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
647 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,
Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay,

- 648 Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson,
649 Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju
650 Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García,
651 Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeon-
652 taek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal,
653 Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor
654 Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou,
655 Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kath-
656 leen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scal-
657 ing language modeling with pathways. *J. Mach. Learn. Res.*, 24:240:1–240:113, 2022. URL
658 <https://api.semanticscholar.org/CorpusID:247951931>.
- 659 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,
660 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned lan-
661 guage models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- 662 Jaeyoung Chung, Suyoung Lee, Hyeongjin Nam, Jaerin Lee, and Kyoung Mu Lee. Luciddreamer:
663 Domain-free generation of 3d gaussian splatting scenes. *arXiv preprint arXiv:2311.13384*, 2023.
- 664 Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig
665 Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of anno-
666 tated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
667 Recognition*, pp. 13142–13153, 2023.
- 668 Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan
669 Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of
670 10m+ 3d objects. *Advances in Neural Information Processing Systems*, 36, 2024.
- 671 Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei,
672 Songyang Zhang, Haodong Duan, Maosong Cao, Wenwei Zhang, Yining Li, Hang Yan, Yang
673 Gao, Xinyue Zhang, Wei Li, Jingwen Li, Kai Chen, Conghui He, Xingcheng Zhang, Yu Qiao,
674 Dahua Lin, and Jiaqi Wang. Internlm-xcomposer2: Mastering free-form text-image composition
675 and comprehension in vision-language large model, 2024.
- 676 Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann,
677 Thomas B McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset
678 of 3d scanned household items. In *2022 International Conference on Robotics and Automation
679 (ICRA)*, pp. 2553–2560. IEEE, 2022.
- 680 Danny Driess, F. Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter,
681 Ayzaan Wahid, Jonathan Tompson, Quan Ho Vuong, Tianhe Yu, Wenlong Huang, Yevgen Cheb-
682 otar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman,
683 Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Peter R. Florence. Palm-e: An
684 embodied multimodal language model. In *International Conference on Machine Learning*, 2023.
685 URL <https://api.semanticscholar.org/CorpusID:257364842>.
- 686 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
687 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for
688 high-resolution image synthesis. *arXiv preprint arXiv:2403.03206*, 2024.
- 689 Daniel Gatis etc. Rembg. <https://github.com/danielgatis/rembg>, 2020.
- 690 Ye Fang, Zeyi Sun, Tong Wu, Jiaqi Wang, Ziwei Liu, Gordon Wetzstein, and Dahua Lin. Make-it-
691 real: Unleashing large multimodal model’s ability for painting 3d objects with realistic materials,
692 2024.
- 693 Junlin Han, Filippos Kokkinos, and Philip Torr. Vfusion3d: Learning scalable 3d generative models
694 from video diffusion models. *arXiv preprint arXiv:2403.12034*, 2024.
- 695 Amir Hertz, Kfir Aberman, and Daniel Cohen-Or. Delta denoising score. In *Proceedings of the
696 IEEE/CVF International Conference on Computer Vision*, pp. 2328–2337, 2023.

- 702 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
703 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*
704 *neural information processing systems*, 30, 2017.
- 705
706 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
707 *neural information processing systems*, 33:6840–6851, 2020.
- 708
709 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
710 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom
711 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia
712 Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and
713 L. Sifre. Training compute-optimal large language models. *ArXiv*, abs/2203.15556, 2022. URL
714 <https://api.semanticscholar.org/CorpusID:247778764>.
- 715 Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli,
716 Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint*
717 *arXiv:2311.04400*, 2023.
- 718
719 Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao
720 Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, Kriti Aggarwal, Zewen Chi, Johan Bjorck,
721 Vishrav Chaudhary, Subhojit Som, Xia Song, and Furu Wei. Language is not all you need:
722 Aligning perception with language models. *ArXiv*, abs/2302.14045, 2023. URL [https://](https://api.semanticscholar.org/CorpusID:257219775)
723 api.semanticscholar.org/CorpusID:257219775.
- 724
725 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan
726 Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi,
727 Ali Farhadi, and Ludwig Schmidt. Openclip. [https://github.com/mlfoundations/](https://github.com/mlfoundations/open_clip)
[open_clip](https://github.com/mlfoundations/open_clip), 2021.
- 728
729 Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint*
730 *arXiv:2305.02463*, 2023.
- 731
732 Yash Kant, Ziyi Wu, Michael Vasilkovsky, Guocheng Qian, Jian Ren, Riza Alp Guler, Bernard
733 Ghanem, Sergey Tulyakov, Igor Gilitschenski, and Aliaksandr Siarohin. Spad: Spatially aware
734 multiview diffusers. *arXiv preprint arXiv:2402.05235*, 2024.
- 735
736 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splat-
737 ting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.
- 738
739 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
740 *arXiv:1312.6114*, 2013.
- 741
742 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
743 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceed-*
744 *ings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.
- 745
746 Jeong-gi Kwak, Erqun Dong, Yuhe Jin, Hanseok Ko, Shweta Mahajan, and Kwang Moo Yi. Vivid-
747 1-to-3: Novel view synthesis with video diffusion models. *arXiv preprint arXiv:2312.01305*,
748 2023.
- 749
750 Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground
751 v2.5: Three insights towards enhancing aesthetic quality in text-to-image generation, 2024.
- 752
753 Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan
754 Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3d: Fast text-to-3d with sparse-view gen-
755 eration and large reconstruction model. *arXiv preprint arXiv:2311.06214*, 2023a.
- 756
757 Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. Blip: Bootstrapping language-
758 image pre-training for unified vision-language understanding and generation. In *International*
759 *Conference on Machine Learning*, 2022. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:246411402)
760 [CorpusID:246411402](https://api.semanticscholar.org/CorpusID:246411402).

- 756 Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. Blip-2: Bootstrapping
757 language-image pre-training with frozen image encoders and large language models. *ArXiv*,
758 abs/2301.12597, 2023b. URL [https://api.semanticscholar.org/CorpusID:
759 256390509](https://api.semanticscholar.org/CorpusID:256390509).
- 760 Yixun Liang, Xin Yang, Jiantao Lin, Haodong Li, Xiaogang Xu, and Yingcong Chen. Lucid-
761 dreamer: Towards high-fidelity text-to-3d generation via interval score matching. *arXiv preprint*
762 *arXiv:2311.11284*, 2023.
- 764 Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten
765 Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d con-
766 tent creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
767 *Recognition*, pp. 300–309, 2023.
- 768 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances*
769 *in neural information processing systems*, 36, 2024a.
- 771 Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen,
772 Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3d objects with
773 consistent multi-view generation and 3d diffusion. *arXiv preprint arXiv:2311.07885*, 2023a.
- 775 Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. One-
776 2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *Advances in*
777 *Neural Information Processing Systems*, 36, 2024b.
- 778 Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick.
779 Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International*
780 *Conference on Computer Vision*, pp. 9298–9309, 2023b.
- 781 Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang.
782 Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint*
783 *arXiv:2309.03453*, 2023c.
- 785 Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma,
786 Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. Wonder3d: Single image to 3d
787 using cross-domain diffusion. *arXiv preprint arXiv:2310.15008*, 2023.
- 788 Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with pre-
789 trained models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 791 Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, Natalia Neverova, Andrea Vedaldi, Oran Gafni,
792 and Filippos Kokkinos. Im-3d: Iterative multiview diffusion and reconstruction for high-quality
793 3d generation. *arXiv preprint arXiv:2402.08682*, 2024.
- 795 Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and
796 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications*
797 *of the ACM*, 65(1):99–106, 2021.
- 798 OpenAI. Gpt-4v(ision) system card. *OpenAI*, 2023a. URL [https://api.
799 semanticscholar.org/CorpusID:263218031](https://api.semanticscholar.org/CorpusID:263218031).
- 800 R OpenAI. Gpt-4 technical report. *arXiv*, pp. 2303–08774, 2023b.
- 802 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
803 *the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- 805 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
806 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
807 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 808 Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d
809 diffusion. *arXiv preprint arXiv:2209.14988*, 2022.

- 810 Zhangyang Qi, Yunhan Yang, Mengchen Zhang, Long Xing, Xiaoyang Wu, Tong Wu, Dahua Lin,
811 Xihui Liu, Jiaqi Wang, and Hengshuang Zhao. Tailor3d: Customized 3d assets editing and gen-
812 eration with dual-side images. *arXiv preprint arXiv:2407.06191*, 2024.
- 813
- 814 Lingteng Qiu, Guanying Chen, Xiaodong Gu, Qi Zuo, Mutian Xu, Yushuang Wu, Weihao Yuan,
815 Zilong Dong, Liefeng Bo, and Xiaoguang Han. Richdreamer: A generalizable normal-depth
816 diffusion model for detail richness in text-to-3d. *arXiv preprint arXiv:2311.16918*, 2023.
- 817 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
818 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
819 models from natural language supervision. In *International conference on machine learning*, pp.
820 8748–8763. PMLR, 2021.
- 821
- 822 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
823 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-
824 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 825 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv
826 preprint arXiv:2202.00512*, 2022.
- 827
- 828 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
829 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
830 open large-scale dataset for training next generation image-text models. *Advances in Neural
831 Information Processing Systems*, 35:25278–25294, 2022.
- 832 Prafull Sharma, Varun Jampani, Yuanzhen Li, Xuhui Jia, Dmitry Lagun, Fredo Durand, William T
833 Freeman, and Mark Matthews. Alchemist: Parametric control of material properties with diffu-
834 sion models. *arXiv preprint arXiv:2312.02970*, 2023.
- 835
- 836 Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen,
837 Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base
838 model. *arXiv preprint arXiv:2310.15110*, 2023a.
- 839 Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view
840 diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023b.
- 841
- 842 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
843 learning using nonequilibrium thermodynamics. In *International conference on machine learn-
844 ing*, pp. 2256–2265. PMLR, 2015.
- 845 Zeyi Sun, Ye Fang, Tong Wu, Pan Zhang, Yuhang Zang, Shu Kong, Yuanjun Xiong, Dahua Lin, and
846 Jiaqi Wang. Alpha-clip: A clip model focusing on wherever you want, 2023.
- 847
- 848 Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative
849 gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.
- 850
- 851 Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm:
852 Large multi-view gaussian model for high-resolution 3d content creation. *arXiv preprint
853 arXiv:2402.05054*, 2024a.
- 854 Shitao Tang, Jiacheng Chen, Dilin Wang, Chengzhou Tang, Fuyang Zhang, Yuchen Fan, Vikas
855 Chandra, Yasutaka Furukawa, and Rakesh Ranjan. Mvdifffusion++: A dense high-resolution
856 multi-view diffusion model for single or sparse-view 3d object reconstruction. *arXiv preprint
857 arXiv:2402.12712*, 2024b.
- 858 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,
859 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly
860 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 861
- 862 Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, Adam Letts, Yangguang Li, Ding
863 Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. Triposr: Fast 3d object reconstruction
from a single image. *arXiv preprint arXiv:2403.02151*, 2024.

- 864 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
865 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Ar-
866 mand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
867 language models. *ArXiv*, abs/2302.13971, 2023. URL <https://api.semanticscholar.org/CorpusID:257219404>.
868
- 869 Vikram Voleti, Chun-Han Yao, Mark Boss, Adam Letts, David Pankratz, Dmitry Tochilkin, Chris-
870 tian Laforte, Robin Rombach, and Varun Jampani. Sv3d: Novel multi-view synthesis and 3d
871 generation from a single image using latent video diffusion. *arXiv preprint arXiv:2403.12008*,
872 2024.
873
- 874 Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jaco-
875 bian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the*
876 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12619–12629, 2023a.
877
- 878 Peng Wang and Yichun Shi. Imagedream: Image-prompt multi-view diffusion for 3d generation.
879 *arXiv preprint arXiv:2312.02201*, 2023.
- 880 Peng Wang, Hao Tan, Sai Bi, Yinghao Xu, Fujun Luan, Kalyan Sunkavalli, Wenping Wang, Zexi-
881 ang Xu, and Kai Zhang. Pf-lrm: Pose-free large reconstruction model for joint pose and shape
882 prediction, 2023b.
883
- 884 Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Pro-
885 lifidreamer: High-fidelity and diverse text-to-3d generation with variational score distillation.
886 *Advances in Neural Information Processing Systems*, 36, 2024a.
- 887 Zhengyi Wang, Yikai Wang, Yifei Chen, Chendong Xiang, Shuo Chen, Dajiang Yu, Chongxuan Li,
888 Hang Su, and Jun Zhu. Crm: Single image to 3d textured mesh with convolutional reconstruction
889 model. *arXiv preprint arXiv:2403.05034*, 2024b.
890
- 891 Zhouxia Wang, Ziyang Yuan, Xintao Wang, Tianshui Chen, Menghan Xia, Ping Luo, and Ying
892 Shan. Motionctrl: A unified and flexible motion controller for video generation. *arXiv preprint*
893 *arXiv:2312.03641*, 2023c.
- 894 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
895 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
896 *neural information processing systems*, 35:24824–24837, 2022.
897
- 898 Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan, Valentin Deschaintre, Kalyan Sunkavalli,
899 Hao Su, and Zexiang Xu. Meshlrm: Large reconstruction model for high-quality mesh, 2024.
900
- 901 Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Jiawei Ren, Liang Pan, Wayne Wu, Lei Yang, Jiaqi
902 Wang, Chen Qian, et al. Omniobject3d: Large-vocabulary 3d object dataset for realistic percep-
903 tion, reconstruction and generation. In *Proceedings of the IEEE/CVF Conference on Computer*
904 *Vision and Pattern Recognition*, pp. 803–814, 2023.
- 905 Tong Wu, Guandao Yang, Zhibing Li, Kai Zhang, Ziwei Liu, Leonidas Guibas, Dahua Lin, and
906 Gordon Wetzstein. Gpt-4v (ision) is a human-aligned evaluator for text-to-3d generation. *arXiv*
907 *preprint arXiv:2401.04092*, 2024.
908
- 909 Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong
910 Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE*
911 *conference on computer vision and pattern recognition*, pp. 1912–1920, 2015.
- 912 Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. Instantmesh:
913 Efficient 3d mesh generation from a single image with sparse-view large reconstruction models.
914 *arXiv preprint arXiv:2404.07191*, 2024a.
915
- 916 Yinghao Xu, Zifan Shi, Wang Yifan, Hansheng Chen, Ceyuan Yang, Sida Peng, Yujun Shen, and
917 Gordon Wetzstein. Grm: Large gaussian reconstruction model for efficient 3d reconstruction and
generation. *arXiv preprint arXiv:2403.14621*, 2024b.

918 Fan Yang, Jianfeng Zhang, Yichun Shi, Bowen Chen, Chenxu Zhang, Huichao Zhang, Xiaofeng
919 Yang, Jiashi Feng, and Guosheng Lin. Magic-boost: Boost 3d generation with mutli-view condi-
920 tioned diffusion. *arXiv preprint arXiv:2404.06429*, 2024.

921

922 Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan.
923 Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In *Proceedings of*
924 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1790–1799, 2020.

925 Biao Zhang, Jiapeng Tang, Matthias Niessner, and Peter Wonka. 3dshape2vecset: A 3d shape
926 representation for neural fields and generative diffusion models. *ACM Transactions on Graphics*
927 *(TOG)*, 42(4):1–16, 2023a.

928

929 Pan Zhang, Xiaoyi Dong, Bin Wang, Yuhang Cao, Chao Xu, Linke Ouyang, Zhiyuan Zhao,
930 Haodong Duan, Songyang Zhang, Shuangrui Ding, Wenwei Zhang, Hang Yan, Xinyue Zhang,
931 Wei Li, Jingwen Li, Kai Chen, Conghui He, Xingcheng Zhang, Yu Qiao, Dahua Lin, and Jiaqi
932 Wang. Internlm-xcomposer: A vision-language large model for advanced text-image comprehen-
933 sion and composition, 2023b.

934 Xin-Yang Zheng, Hao Pan, Yu-Xiao Guo, Xin Tong, and Yang Liu. Mvd²: Efficient multiview 3d
935 reconstruction for multiview diffusion. *arXiv preprint arXiv:2402.14253*, 2024.

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

A APPENDIX

A.1 EVALUATION ON WILD PROMPTS FROM REAL USERS

The results of the main part of the paper are only tested on GPT generated prompts. To test our work's capability in wild cases, we also collect real user prompts and compare our method with Instant3D (Li et al., 2023a). Specifically, we randomly collect 100 prompts from <https://www.meshy.ai/> and test the CLIP-R precision as well as GPT based evaluation (detailed in Sup. A.6). Results and some qualitative cases are shown in Tab. 4 and Fig. 8. We highlight that our Bootstrap3D excels Instant3D (Li et al., 2023a) when tested on real user prompts through training on synthetic data.

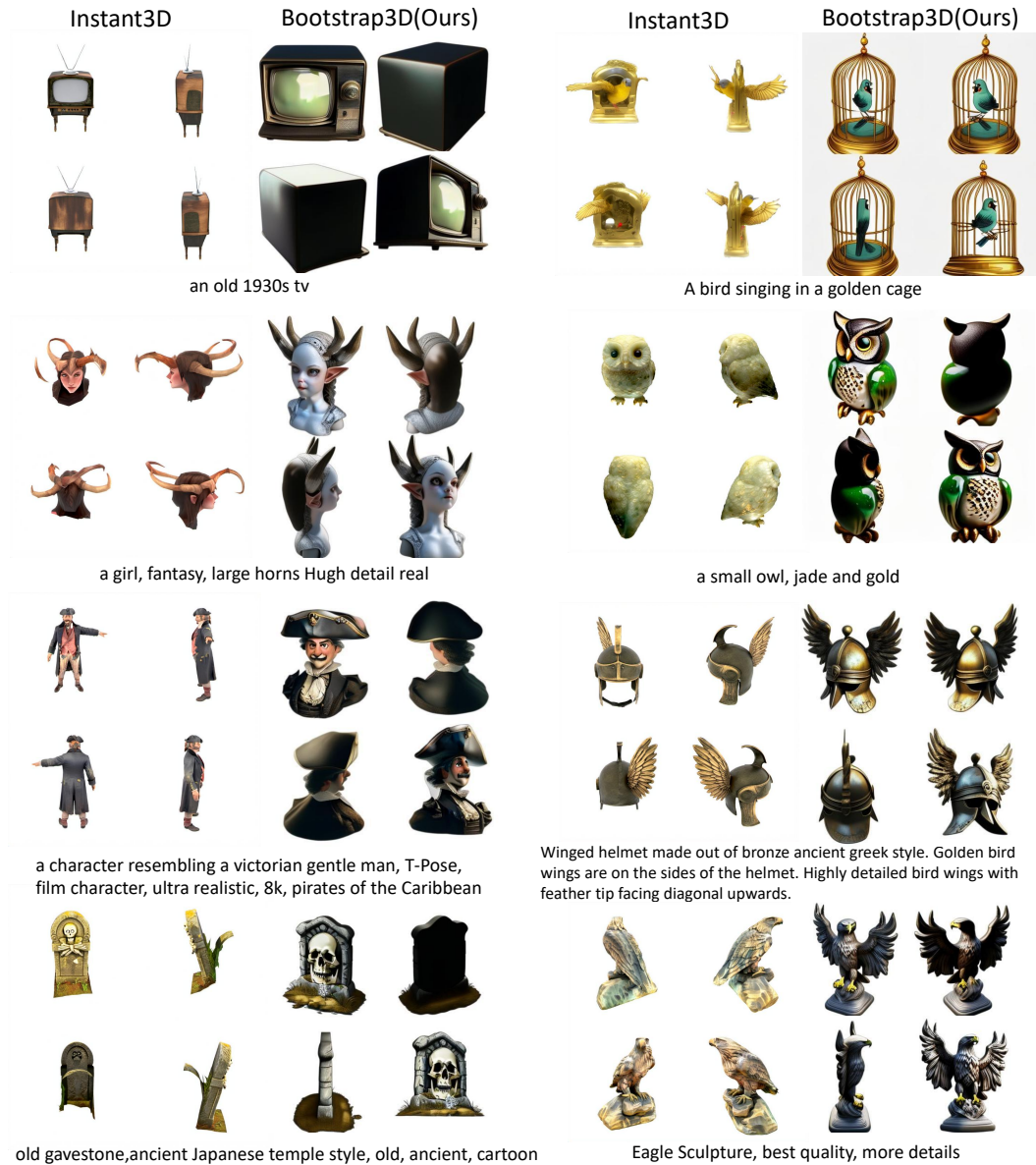


Figure 8: **Real user prompt cases** visualization compared to Instant3D (Li et al., 2023a)

A.2 MORE VISUALIZATION COMPARED TO OTHER METHODS.

We show more visualization of the quantitative experiments shown in the main paper in Fig.9

Table 4: **Test results of in the wild cases.** Bootstrap3D also excels Instant3D (Li et al., 2023a) in generating high quality images according to real user prompts.

Method	CLIP based metric	GPTEval3D	
	CLIP-R score	image-text alignment	texture detail
Instant3D (unofficial)	77.0	22.0%	24.5%
Bootstrap3D	83.5	78.0%	75.5%

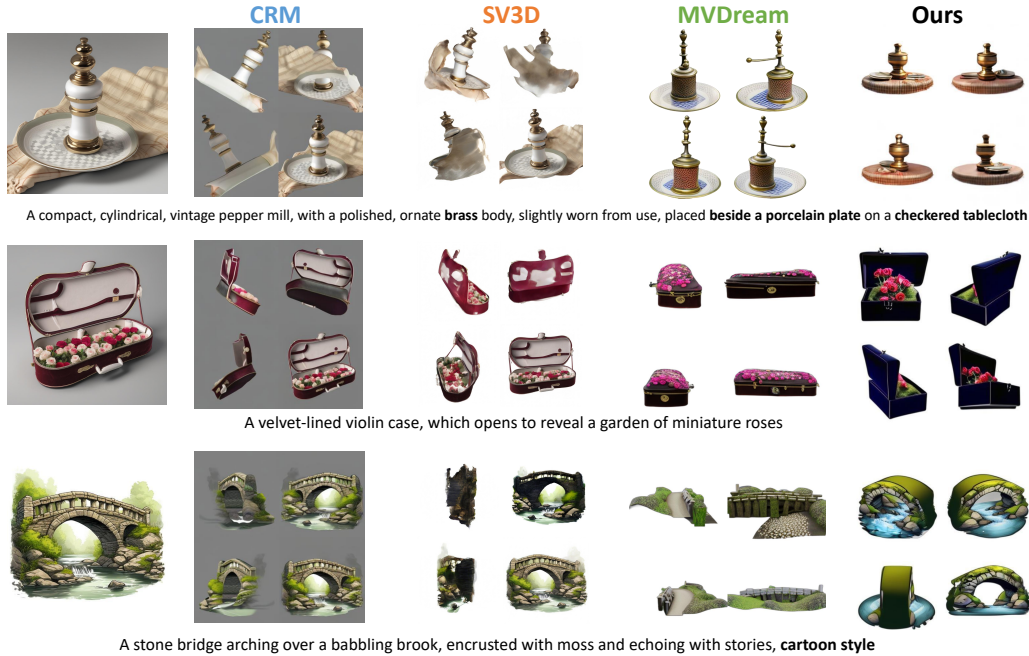


Figure 9: **Generated multiview images compare to other methods.** Our method can generate multi-view images with long text control without encountering blurring effect from data generated by SV3D thanks to TTR and quality filtering.

For Image-to-3D methods, they can sometimes produces significant motion blurring and fails when the input image is out-of-distribution (like the 3rd cartoon style case). We resample the high-quality segment of the distribution of generated images using quality filtering based on MLLM methods. Furthermore, by employing TTR, we limit the impact of these data when training multi-view diffusion models, allowing our model to produce much clear results. In addition, we use a caption rewriting method, enabling finer prompt control for the generated multi-view images.

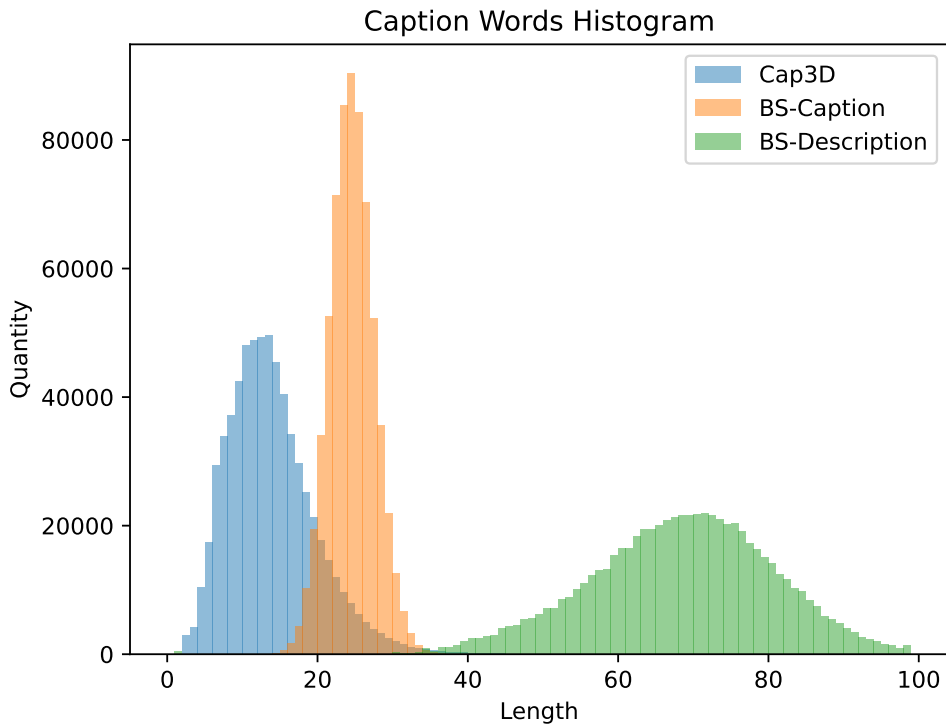
A.3 DATA STATISTICS

A.3.1 CAPTION ANALYSIS

Fig. 10 and 11 provide a visualization of the root noun-verb pairs for the captions generated by GPT-4V (OpenAI, 2023a) and MV-LLaVA. It’s clear to see that the diversity and linguistic expression of the captions produced by MV-LLaVA are highly matched with those of GPT-4V. We believe the highly detailed description focusing on object’s texture, shape and color have potential usage beyond training multi-view diffusion model in the field like object texturing (Fang et al., 2024) and stylization (Sharma et al., 2023) in Computer Graphics. MV-LLaVA can also serve as free and efficient 3D object assistant comparable with GPT-4V for future research of 3D content creation.

Fig. 12 visualizes the histogram of caption length compared with Cap3D (Luo et al., 2024). We fine-tune MV-LLaVA to generate two different lengths suitable for different diffusion architecture, namely CLIP-based text encoding (Blattmann et al., 2023; Podell et al., 2023) with 77 token length and T5 based text encoding (Chen et al., 2023a; 2024a) with 120 token length. Both excel the length of Cap3D with less hallucinations.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159



1160 Figure 12: **Histogram Visualization of the Caption Length** compared with Cap3D (Luo et al.,
1161 2024)

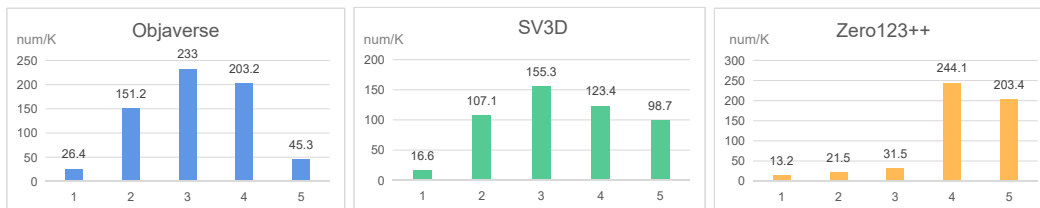
1163 Table 5: **Comparison of lexical composition of the captions** generated by GPT4-Vision and Share-
1164 Captioner.

Lexical	n.	adj.	adv.	v.	num.	prep.
GPT-4V (OpenAI, 2023a)	29.1%	16.0%	1.5%	11.1%	0.5%	9.0%
BS-Description	28.5%	16.0%	1.4%	10.8%	0.3%	8.6%
BS-Caption	30.2%	23.0%	0.3%	5.6%	0.1%	8.9%

1165
1166
1167
1168
1169
1170
1171

1172 Fig. 16, Fig. 17 and Fig. 18. We use high quality data with score 4 and 5 for the training of multi-
1173 view diffusion model.

1174
1175
1176
1177
1178
1179
1180
1181



1182 Figure 13: **Quality score statistics of different data source.**

1183
1184
1185
1186
1187

We count the number of multi-view images from different data sources, namely 660K from Objaverse, 500K from SV3D (Voleti et al., 2024) and 500K from Zero123++ (Shi et al., 2023a) generated by our Bootstrap3D pipeline. Result are shown in Fig.13. For Objaverse and SV3D, the assigned score is normal and we use score 4 and score 5 multi-view images as high quality data for training.

However, for Zero123++, most objects are assigned with score greater than 3. We attribute this phenomenon to the fact that Zero123++ tend to generate objects with less motion blurring but more stretching and deformation compared to SV3D. Joint training of MV-LLaVA on three different data source lead to higher and more focused distribution for Zero123++’s multi-view images. For this part of synthetic data, we leave only score 5 multi-view images as high quality data.

A.4 QUALITY OF MV-LLAVA

A.4.1 CHOICE OF NUMBER OF UNFROZEN LAYERS OF VISION ENCODER.




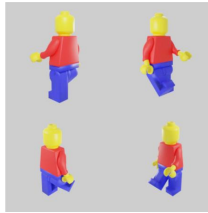
input image				
unfreeze 0 layers	The model displays a basic medieval helmet in matte olivegreen with a visor pattern and nasal guard, suitable for simple visualizations or game assets with a fantasy style. Tag: [single object] [Cartoon]	A cartoonish turtle model with a smooth, creamy yellow body, large eyes, and a friendly expression, featuring soft shading and a glossy finish for a playful appearance. Tag: [single object] [Cartoon]	A detailed classic wingback chair model, with realistic dark charcoal fabric and weathered brown wooden legs and frame conveying an authentic and elegant aesthetic. Tag: [single object] [Photo realistic]	Stylized humanoid figure with a playful design, featuring a yellow head and hands, red torso, blue legs, and a smiling face , with a glossy finish for animation or gaming. Tag: [single object] [Cartoon]
unfreeze 8 layers	Stylized humanoid head with a green, matte finish, featuring a white symbol and purple lines, with darker green hair, suitable for fantasy or historical themes. Tag: [single object] [Cartoon]	The model displays a minimalist, cartoonish humanoid with a uniform cream color and simple black eyes, suggesting a basic prototype or abstract character design. Tag: [single object] [Cartoon]	The model showcases a wingback chair with detailed worn leather in dark brown , contrasting matte wooden legs, and a shiny, curved wooden frame, suitable for realistic interior visualizations. Tag: [single object] [Photo realistic]	A stylized humanoid figure with a glossy yellow head, red torso, and blue legs, featuring a minimalist face and smooth surfaces, ideal for animation or game usage. Tag: [single object] [Cartoon]

Figure 14: **Qualitative results of unfreeze final layers of CLIP (Radford et al., 2021) vision encoder** compared to original fixed vision encoder setting in LLaVA (Liu et al., 2024a).

Inspired by ShareGPT-4V (Chen et al., 2023b), we unfreeze selected final layers of the CLIP (Liu et al., 2024a) vision encoder during the initial phase of vision language alignment. The CLIP-L/14 model used for LLaVA (Liu et al., 2024a) contains 24 transformer layers. We selectively unfreeze some of final layers to enable the CLIP model to focus more on details such as texture of multi-view images. After qualitative manual screening, we select to unfreeze eight layers to yield better results. Fig. 14 illustrates the differences between unfreezing eight layers and not unfreezing any (the original training setting of LLaVA (Liu et al., 2024a)). The red sections highlight the erroneous hallucinations occurring when the vision encoder remains fully unchanged, while the green sections indicate accurate descriptions of the image content. This demonstrates that partially unfreezing the vision encoder can produce more precise captions and reduce some hallucinations.

A.4.2 QUANTITATIVE QUALITY STUDY

To test the quality of our MV-LLaVA. We propose two quantitative study over the quality of captions and the alignment of quality estimation with human experts. In first study, we randomly picked 200 object from Objaverse (Deitke et al., 2023) and exclude training data of MV-LLaVA. We use GPT4-V (OpenAI, 2023a) and MV-LLaVA to generate descriptive captions for each object. We invite human volunteers to choose their preference over shuffled captions. Results are shown in Tab. 6, where MV-LLaVA shows comparable captioning ability with powerful GPT4-V (OpenAI, 2023a), which is essential to generate millions of high quality image-text pairs for the training of text to multi-view image diffusion model.

Second experiment studies MV-LLaVA’s ability in quality estimation of both 3D assets and generated multi-view images. We invite human volunteers to estimate the quality of multi-view images rendered from Objaverse (Deitke et al., 2023) or generated by SV3D (Voleti et al., 2024). As there is no golden standard for multi quality classification, We ask them to separate the randomly select

multi-view images into approximately two half and serve as GT quality. We use MV-LLaVA to estimate the quality of these images and generate confusion matrix. Results are shown in Tab.7. Given the great amount of source data of 3D assets and infinite synthetic data, we care more about the false positive rate, as these data will be mixed into training data. In this observation, we highlight the false positive rate of over 20% for SV3D (Voleti et al., 2024) generated multi-view images. This result align with the observation of inevitable motion blurring of SV3D (Voleti et al., 2024). To leverage this part of data source for data diversity without hurting the final quality. We propose Training Noise Reschedule to avoid samplings from these synthetic data when time step is small.

Table 6: **Human evaluation** on the quality of generated captions from MV-LLaVA vs. GPT4-Vision (OpenAI, 2023a) over 200 validation samples from Objaverse (Deitke et al., 2023).

Preference	GPT4-Vision (OpenAI, 2023a)	MV-LLaVA	Comparable
Percentage	39.5%	34.5%	26.0%

Table 7: **Confusion matrix** of mutli-view images quality estimation.

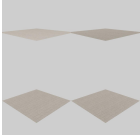


	Objaverse quality check		Synthetic quality check		
	HQ-gt	LQ-gt		HQ-gt	LQ-gt
HQ by model	31.0%	4.5%	HQ by model	34.5%	11.5%
LQ by model	11.0%	53.5%	LQ by model	17.0%	37.0%

A.4.3 QUALITATIVE CAPTION QUALITY STUDY




We selective compare some of the captions generated by Cap3D (Luo et al., 2024) and MV-LLaVA in Fig. 15. Our MV-LLaVA can generate more detailed descriptive captions with less hallucinations.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

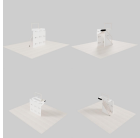

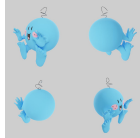
Score: 1

 <p>Cap3D: 3D rendering of a square, tiled floor and ceiling.</p>	<p>Ours: The model shows a plain with nearly identical flat surfaces with a uniform beige color, lacking any distinctive features or textures. There are no discernible details indicating what the surfaces represent or their potential use, and thus the model serves as an extremely basic representation of a flat object.</p>	 <p>Cap3D: A collection of drawings featuring a jellyfish, a girl with various accessories, a fish, and a unicorn.</p>	<p>Ours: The model is a highly abstract representation, possibly of a figure or object with an unclear purpose. It features thin, rod-like lines and sparse clusters of pink and blue colors, lacking any discernible texture, material quality, or detailed features.</p>	 <p>Cap3D: A 3D object featuring elements of maps, a fish, a heart, a person, and a ghost, all adorned with pink and blue dots.</p>	<p>Ours: The model is composed of faint, scattered dots and lacks discernible shape or texture, making it difficult to determine its intended representation. It appears as a collection of sparse, flat details without any clear function or recognizability.</p>
---	--	--	---	--	--




Score: 2

 <p>Cap3D: A 3D rendering of a white exhibition space featuring a round table, pictures, a water feature, and a chair.</p>	<p>Ours: The model represents a minimalist art gallery scene with discernible but simplified elements. Semi-transparent white walls and floors suggest a clean, empty space. Slight color variations hint at lighting nuances. Sparse objects including picture frames are rendered in simplistic grayscale, lacking texture but offering a rudimentary representation of an art display.</p>	 <p>Cap3D: A low poly 3D model of a white ghost, small bird, and teddy bear.</p>	<p>Ours: The model is a low-polygon representation of a character's head with a simplified form, featuring a broad, flat surface with minimal texture. The main body of the model is white, while the facial features are depicted in a flat grey tone, lacking depth and detail. It suggests a basic, monochromatic material, possibly indicative of clay or a digital sculpting tool...</p>	 <p>Cap3D: 3D models of a fruit, paper airplane, rock, and green leaf with various surfaces and interactions.</p>	<p>Ours: The model represents a rock-like object with a textured surface, depicting moss or lichen patches in shades of green and gray. The underlying stone surface exhibits rust tones suggesting iron-rich sedimentary rock. However, the model lacks finer details such as specific rock strata or additional variations that could enhance its realism.</p>
--	--	--	--	--	---

Score: 3

 <p>Cap3D: A 3D rendering of a white cabinet, staircase, small box, chair, desk, ceiling-mounted light fixture, and ceiling panel.</p>	<p>Ours: The 3D model represents a minimalist architectural form, possibly a building or pavilion with an external staircase. It's composed of white and grey tones, where white suggests a smooth concrete or plaster and grey represents metallic or stone features. The model lacks intricate textures but has discernible geometric shapes and shadows, offering a basic sense of depth and spatial relationships.</p>	 <p>Cap3D: A 3D wooden toy car model featuring a stove, chair, pipe, spout, bottle, and steering wheel.</p>	<p>Ours: The model represents a whimsical, simplified vehicle consisting of a brown wagon with yellow wheels and a red, toy-like trailer being pulled. The wagon's body has a flat, wood-like texture, while the trailer's body shows a smooth, glossy finish. The handle and wheels of the trailer are colored in bright yellow, enhancing its toy-like appearance...</p>	 <p>Cap3D: A 3D blue sphere with pink eyes and hands, resembling a cartoon character.</p>	<p>Ours: The model is a whimsical, cartoon-style character with a spherical blue body and exaggerated, bulbous limbs. The facial features include a wide, cheerful smile and large, expressive eyes with a simple white sclera and black pupils. The limbs' ends have been stylized into points, enhancing the playful nature of the model...</p>
---	---	--	---	---	--

Score: 4

 <p>Cap3D: 3D model of a Japanese garden village with trees, a bridge, and a hill, featuring a plane flying overhead.</p>	<p>Ours: This model portrays a small, stylized scene of a Japanese garden with a pond and cherry blossom trees. The pond is depicted in a tranquil shade of blue, surrounded by vibrant green grass. Pink cherry blossom trees, a symbol of spring in Japan, are scattered around the pond with varying shades of pink. A small bridge with a red railing crosses the pond, adding a pop of color to the scene. The textures are simple yet distinct, with the grass and trees having a flat appearance...</p>	 <p>Cap3D: A 3D model of a white, yellow, and gold machine-like spaceship.</p>	<p>Ours: This model is a robotic chair with a minimalist design, rendered in a clean, modern style. The seat and backrest are white, implying a smooth, possibly plastic or painted metal surface. The mechanical arm and base are a metallic gray, suggesting a metallic texture with a matte finish. Gold-colored accents, likely meant to represent controls or adjustments, add a touch of color to the model...</p>	 <p>Cap3D: Yvy toilet paper box, Yui air purifier, Yoyo toner cartridges, Yui vitamins, and Yui sanitary napkins in yellow and white packaging with respective logos.</p>	<p>Ours: The model is a packaged product, most likely an air purifier, with a recognizable brand and model identifiable from the images. The packaging is a vibrant yellow with white accents, which gives a sense of a plastic or laminated paper material. The product image on the box shows a white device with a blue air filter, conveying a clean and modern aesthetic...</p>
---	---	--	---	--	---

Score: 5




 <p>Cap3D: A 3D horn-shaped jug with a wooden handle and writing on it.</p>	<p>Ours: This 3D model represents a wooden tankard with a detailed wood grain texture in natural tan shades, featuring realistic knots and variations in the wood. The handle and rim of the tankard are well-defined, and the model exhibits a rustic, handcrafted appearance. The carved lines and wood texture give it an authentic look...</p>	 <p>Cap3D: Red Peugeot 208 SUV 3D model.</p>	<p>Ours: The model is a red SUV with a high level of detail, including realistic shapes for the body, windows, and wheels. The body has a glossy finish, simulating automotive paint, while the black trim and silver rims have a matte texture. The wheels are detailed with intricate silver spokes that resemble alloy wheels, and the tires are a vibrant green with tread patterns adding to the model's realism...</p>	 <p>Cap3D: 3D model of a wagon with a roof and wheels, featuring a horse.</p>	<p>Ours: This 3D model is a charming, medieval-style vendor cart with a high degree of texture detail, suggesting materials like weathered wood for the main structure, rusty metal for the fittings and wheel rims, and tattered fabric for the awning. The cart is laden with goods such as sacks, barrels, and hanging items...</p>
---	---	--	---	--	---

Figure 15: Caption comparison with Cap3D (Luo et al., 2024). Our MV-LLaVA can generate long captions that faithfully describing 3D assets from different perspectives like color, geometry and texture.



Figure 16: Randomly picked multi-view images with different scores from 500k synthetic data generated by SV3D (Voleti et al., 2024).

1403



Figure 17: Randomly picked multi-view images with different scores from 500k synthetic data generated by Zero123++ (Shi et al., 2023a).

1455

1456

1457



Figure 18: Randomly picked multi-view images with different scores from 660k Objaverse (Deitke et al., 2023) 3D assets.

1511

A.5 DETAILS OF PROMPT DESIGN

A.5.1 PROMPTS FOR GPT-4V FOR QUALITY CHECK

Assume you are a quality checker of a diffusion model. This diffusion model is trained to achieve novel view synthesis. I give this model the image in the upper-left side and it generate novel views in the rest three images(upper-right, lower-left, lower-right). You should tell me the quality of the generated novel view images. The score ranges from 1 to 5, representing the quality of the model from low to high. The detailed evaluation criteria are as follows:

1. The novel views are difficult to discern what the image supposed to be, lacking in recognizability. It has no usable value.
2. The novel views are distinguishable, clearly determine what the object/scene is similar to the given ground truth image. However, there is obvious inconsistency between the novel view synthesized images and ground-truth image. There are many obvious areas of image is blurred or indicating rotation.
3. The novel views are relatively good, the inconsistency between novel view synthesized images with ground-truth image is not obvious. The blurring area indicating rotation or uncertainty is acceptable for usage.
4. The novel views are pretty good, although the might be blurring areas or less resolution, the view consistency is well maintained.
5. The novel views are excellent. It is hard to tell which image from four is ground-truth and which is synthesised.

You should give me the overall score with one score number, with reason in next line. besides the quality check, I need you to generate a long descriptive caption for the scene/object from 4 different view, focusing on the part/object relative position, color, number of objects and so on with no more than 50 words and no less than 30 words. DO NOT MENTION MULTI-VIEW IMAGES FROM DIFFERENT PERSPECTIVE since it is a single scene/object. you should rearrange your result in a JSON format. if all the images(include the ground-truth image) are of low quality, just output a lowest score.

Here is an example for you:

```
{
  "score": 4,
  "reason": "The novel views generated from the model are quite convincing with a high degree of consistency in terms of texture, lighting, and color when compared to the ground-truth image. There is some minor distortion in shape and perspective, but the overall quality is high, and it maintains the realism of the scene.",
  "caption": "A cluster of shiny five apples, ranging from deep red to sunny yellow, sits comfortably within a rustic woven basket. Their smooth, round forms are grouped closely, reflecting light and casting soft shadows that accentuate their voluminous curves and vibrant colors."
}
```

This is a quad image generated from rendering a SINGLE 3D model FROM FOUR DIFFERENT views. I would like you to score the quality of this models to evaluate its current state. The score ranges from 1 to 5, representing the quality of the model from low to high. The detailed evaluation criteria are as follows:

- 1 point: The overall quality of the model is quite poor, making it difficult to discern what it is supposed to be, lacking in recognizability. The model is almost one solid block, or extremely scattered, or in fragments. It has no usable value.
- 2 points: The overall quality of the model is relatively poor, but it is possible to guess what it is, possessing low recognizability. It preliminarily has some geometric shape and can be considered a prototype model element, lacking identifiable material information, and almost has no usable value.
- 3 points: The overall quality of the model is average, it is possible to determine what it is, having certain recognizability. Different areas use different materials (colors), it preliminarily has usable value, and initially has aesthetic value.
- 4 points: The overall quality of the model is relatively high, it can be clearly determined what it is, with high recognizability. It preliminarily has certain texture details, and different parts of a model can be clearly distinguished, having high usable value and certain aesthetic value.
- 5 points: The overall quality of the model is extremely high, allowing for the classification of the model's type at a very fine granularity. It has high texture details, is a fully formed 3D model that can be used for games, simulations, or even animations, and has high aesthetic value.

After scoring, please also generate a description of the current model. If the model quality is low, only a brief description is needed; when the model quality is high, a complete description of the different details of the model is required. The description process should focus on color, material, texture details as much as possible. You can also recommended to suggest overall style. With NO MORE THAN 120 words. Especially describe color and material of different parts concretely and faithfully, let the reader easily imagine the same model.

Finally, I hope you can annotate two kinds of tags for the model. Tag1 is about the style of overall model. You can choose from [photo-realistic], [carton] and [CAD]. Tag model as [CAD] when seems like a preliminary work build by CAD software and not real. Tag model [carton] when it is good enough with carton style. Tag model [photo-realistic] when model seems like real object in the world; Tag2 is about the scale that the model represents, you can choose from [single object], [multi-object], [small scene] and [large scene]. Assign model [large scene] when it represents scene like urban street, park, etc. Assign it as [small scene] when it represents scene like inner structure or design of a house, small area, etc. Assign it as [multi-object] when it represents combination of multi objects. Assign it as [single object] when it represents single object.

Here are three examples. You should follow this format:

e.g. 1
Score: 1
Description: The model depicts a very basic and abstract urban planning concept with indistinct structures and simplistic landscaping, lacking detail and texture, appropriate for early-stage design or conceptual visualization.
Tag: [Photorealistic] [large scene]

e.g. 2
Score: 2
Description: The object is a simple sphere with a homogeneous speckled texture, suggesting a stone-like material. The colors vary slightly between shades of dark gray, brown, and rust, with a matte finish. It lacks specific features or details that would indicate a higher level of complexity or function.
Tag: [Photorealistic] [single object]

e.g. 3
Score: 3
Description: The model appears to represent an architectural structure with two levels. Different colors suggest varied materials: translucent white for the structural framework, solid blue representing walls or glass panels, and yellow for interior elements, possibly stairs or floors. The style seems utilitarian, potentially for preliminary construction visualization.
Tag: [CAD] [small scene]

e.g. 4
Score: 4
The model depicts a metallic livestock handling equipment known as a cattle chute. It is rendered in shades of dark gray, conveying a metallic texture consistent with steel or iron. The structure is detailed with bolts, bars, and sliding gates, implying a sturdy construction. Text labels like "METALCORP" and "CATTLE MASTER!" in blue enhance realism, suggesting a commercial quality model suitable for simulations or instructional material. The style is industrial and pragmatic.
Tag: [Photorealistic] [single object]

e.g. 5
Score: 5
The model is a stylized, anime-inspired character with a cheerful expression. Hair is rendered in a turquoise shade, contrasting with ribbons in alternate hues of pink and blue. Skin tone is in a soft peach, while the outfit combines white, grey, and gold tones, with a large yellow flower accessory. Surfaces show subtle shading, indicating variations in material. The playful, colorful appearance suggests a light-hearted, fantasy aesthetic.
Tag: [Cartoon] [single object]

Figure 19: Prompt for GPT-4V to generate caption and estimate quality of multi-view images from SV3D (Voleti et al., 2024), zero123++ (Shi et al., 2023a) and Objaverse (Deitke et al., 2023).

Detailed prompts are shown in Fig.19.

1566 A.5.2 PROMPTS FOR MV-LLAVA INSTRUCT TUNING
1567
15681569 Table 8: **Instruct tuning prompt for SV3D (Voleti et al., 2024) and Zero123++ (Shi et al., 2023a)**
1570 **multi-view images**

prompt type	prompt
generate caption	<image><image><image><image>\nWhat is this multi-view photo about? generate a short caption for me. <image><image><image><image>\nGenerate a short caption of the following multi-view image. <image><image><image><image>\nCan you describe the main features of this multi-view image for me by a short caption?
reasoning	How about the view consistency of this synthesized multi-view image? Do some comments about the view consistency of this synthesized multi-view image. What do you think about the view consistency of this synthesized multi-view image?
quality estimation	What do you think about the overall quality of view consistency of three synthesized novel views? Choosing from "poor", "relatively poor", "boardline", "relatively good", "good", "perfect".

1587
1588 Table 9: **Instruct tuning prompt for Objaverse (Deitke et al., 2023) rendered multi-view images**
1589

prompt type	prompt
long description	<image><image><image><image>\nWhat is this multi-view photo about? generate a long descriptive caption for me. <image><image><image><image>\nGenerate a long descriptive caption of the following multi-view image. <image><image><image><image>\nCan you describe the main features of this multi-view image for me by a long descriptive caption?
caption	<image><image><image><image>\nWhat is this multi-view photo about? generate a short caption for me. <image><image><image><image>\nGenerate a short caption of the following multi-view image. <image><image><image><image>\nCan you describe the main features of this multi-view image for me by a short caption?
quality estimation	What do you think about the overall quality of this 3D model? Choosing from "poor", "relatively poor", "boardline", "relatively good", "good", "perfect".
scale tag	What do you think about the scale of the 3D model represents? Choosing from "single_object", "multi-object", "small_scene", "large_scene".
style tag	What do you think about the overall style of the 3D model? Choosing from "CAD", "Cartoon", "Photo_realistic".

1612
1613 A.6 GPT-4V BASED 3D OBJECT GENERATION EVALUATION.
1614

1615 We adopt method proposed in GPTeval3D (Wu et al., 2024) for more thorough and human-aligned
 1616 evaluation of the quality of generated object by different methods. A full test case is shown in
 1617 Fig. 20. Left 9-view image is rendered from object generated by Bootstrap3D and the right one
 1618 generated by Instant3D (Li et al., 2023a). We ask GPT-4V to mainly evaluate through comparison
 1619 based on three dimensions: text-image alignment, low-level texture quality and 3D plausibility.
 The answer of GPT-4V shows its in depth perception ability of given reasonable comparison well

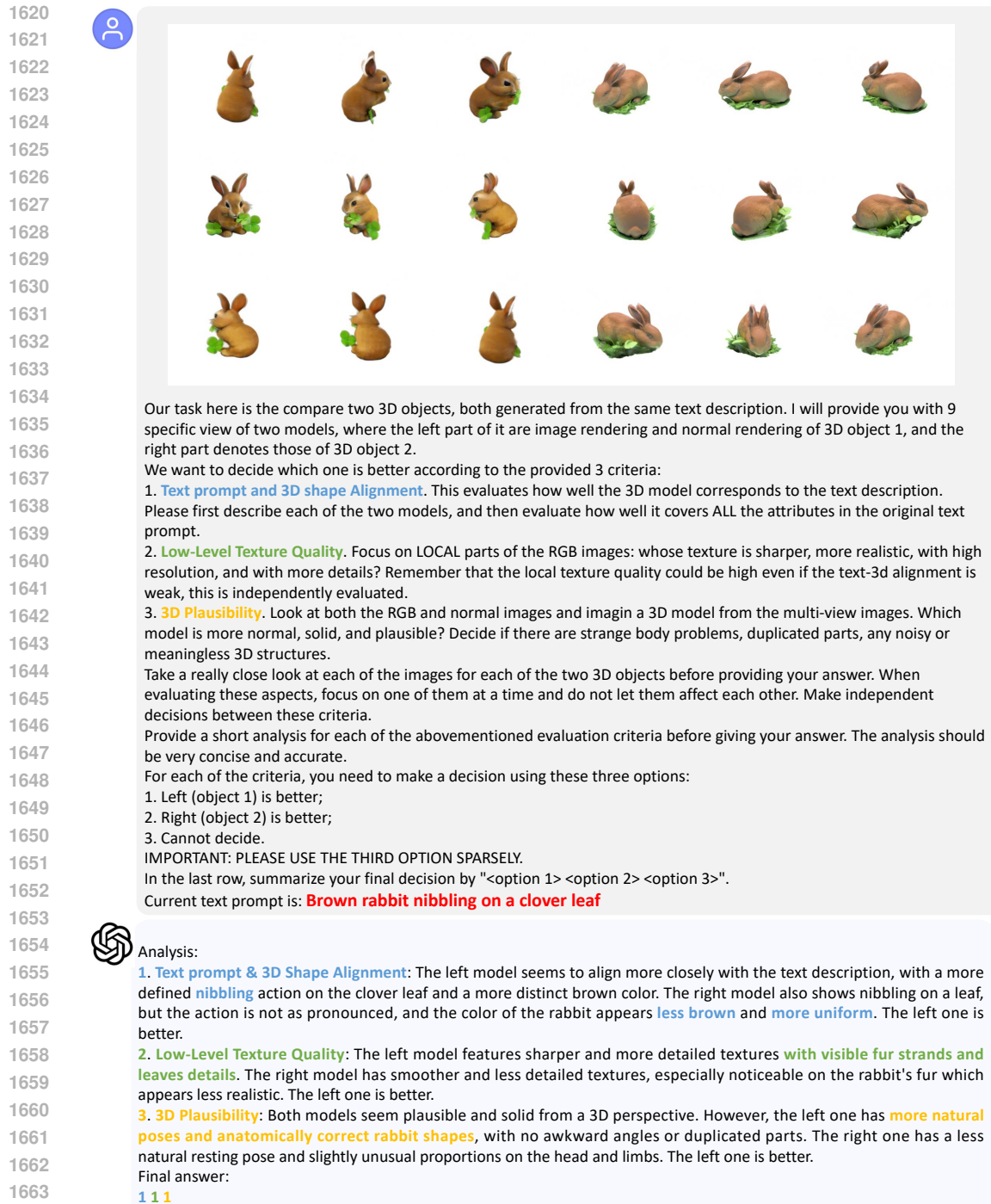


Figure 20: A test conversation with GPT-4V (OpenAI, 2023a) of evaluating generated objects.

aligned with human preference. We thus choose to use GPT-4V rather than human volunteers to give reasonable evaluation.

We adopt the 110 test prompts proposed in GPTeval3D (Wu et al., 2024) to test Bootstrap3D generated object comparing with Instant3D (Li et al., 2023a), Zero123++ (Shi et al., 2023a) and MV-Dream (Shi et al., 2023b). For each methods, we conditioned model based on 110 test prompts with

1674 4 different seeds, with each methods generates 440 objects, we make 1-to-1 comparison following
 1675 aforementioned test setting. Results are reported in Tab. 10. Except MVDream (Shi et al., 2023b)
 1676 (SDS) (which generates single object consuming 30 mins while Bootstrap3D only need 5 seconds.).
 1677 Bootstrap3D excels in all three evaluation dimensions, which proves the ability of Bootstrap3D in
 1678 creating high quality 3D objects.

1679
1680

1681 Table 10: **GPT-4V based evaluation result.** the result is in format of "number of objects preferred
 1682 generated by Bootstrap3D/ that of other methods". Cases when GPT cannot answer the question or
 1683 generates "cannot decide" answer are excluded.

1684

	Image-text alignment	Texture quality	3D plausibility
1685 Compared to Instant3D (Li et al., 2023a) (unofficial)	247 / 116	202 / 162	259 / 110
1686 Compared to Zero123++ (Shi et al., 2023a)	192 / 143	210 / 161	231 / 139
1687 Compared to MVDream (Shi et al., 2023b) (GRM)	290 / 71	245 / 131	284 / 102
1688 Compared to MVDream (Shi et al., 2023b) (SDS)	188 / 155	173 / 190	192 / 150

1689

1690
1691

1692 A.7 IMPROVING DIRECT 3D GENERATIVE MODELS

1693



1702
1703 Figure 21: **Fine tuned Shape-E generation results** that shows better object-text alignment than
 1704 original Shape-E (Jun & Nichol, 2023) and finetuned version in Cap3D (Luo et al., 2024).

1705

1706

1707
1708 Table 11: **Test results on Shape-E.** More accurate and descriptive 3D caption help model to achieve
 1709 better object-text alignment.

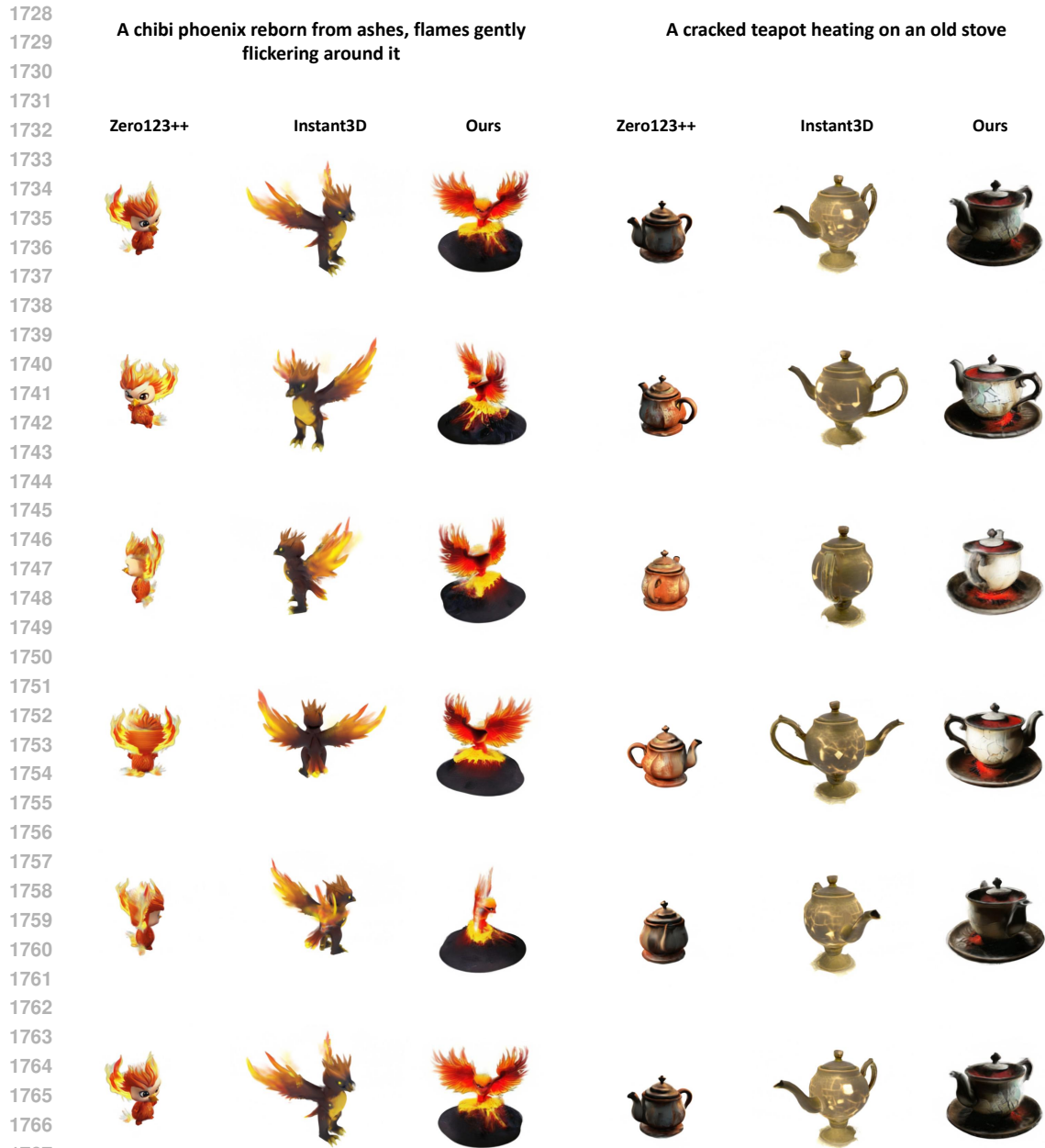
1710

Method	FID ↓	CLIP score ↑	CLIP-R-precision ↑
1711 Shape-E	37.2	80.4	20.3
1712 Cap3D	35.5	79.1	20.0
1713 Ours	35.3	81.2	22.1

1714

1715

1716
1717 In addition to fine-tuning the multiview diffusion model, we also evaluate our framework on direct
 1718 3D generative models, circumventing the use of multi-view images as intermediaries. For this pur-
 1719 pose, we selected the Shape-E (Jun & Nichol, 2023) model for experiment and assess the outcomes
 1720 following the testing method the same to Cap3D (Luo et al., 2024). Specifically, we fine-tune Shape-
 1721 E using 250K BS-Objaverse data, ensuring that all entries scored greater than 3, accompanied by
 1722 more precise and descriptive captions. The metrics for training and testing are consistent with those
 1723 employed in Cap3D (Luo et al., 2024). Some qualitative results are presented in Fig.21, where our
 1724 finetuned version can generate object that follow text prompt more precisely. Quantitative results
 1725 are detailed in Tab.11, where more accurate and descriptive captions than Cap3D can significantly
 1726 improve metrics like CLIP score. Our findings indicate that improved data quality can significantly
 1727 enhance object-text alignment and visual quality of Shape-E. This experiment substantiates that
 our pipeline, characterized by detailed captions and quality filtering, is also effective for direct 3D
 objects generation represented by neural field.



1769 **Figure 22: Visualization of generated objects compared to other edge-cutting methods**

1770

1771

1772 **A.8 MORE RESULTS VISUALIZATION**

1773

1774 **A.8.1 COMPARISON WITH OTHER METHODS**

1775

1776 **A.8.2 VISUALIZATION OF GENERATED OBJECTS WITH DIFFERENT STYLES**

1777

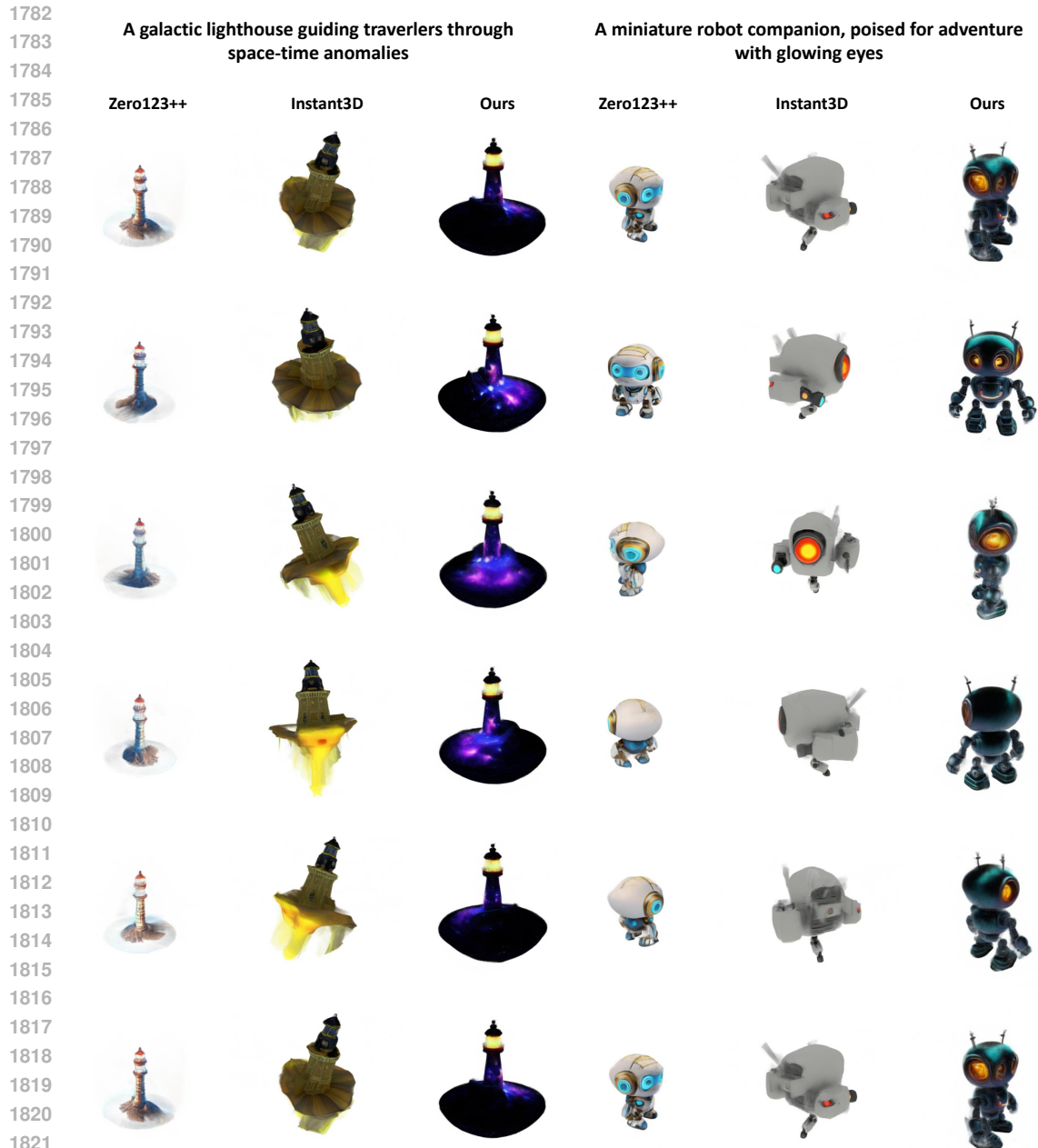
1778 **A.9 BROADER IMPACTS**

1779

1779 **Potential positive societal impacts:** The proposed framework, Bootstrap3D, enhances the quality and consistency of 3D models, which can benefit various industries such as entertainment, education, virtual reality, and digital art. By generating and sharing a large synthetic dataset of high-quality synthetic multi-view images, We will promotes open access to resources that can accelerate progress

1780

1781



1822
1823 **Figure 23: Visualization of generated objects compared to other edge-cutting methods**
1824

1825
1826 in the field. The model and data can serve as educational tools for students and researchers, fostering
1827 learning and innovation in machine learning and 3D modeling.

1828 **Potential negative societal impacts:** High-quality 3D models could be used to create deepfakes or
1829 misleading content, which may contribute to disinformation or malicious activities. Monitoring and
1830 Defense Mechanisms: Developing tools to detect and prevent the misuse of the generated 3D mod-
1831 els, particularly in contexts like disinformation and surveillance. There may be unintended biases
1832 in the generated data or models, leading to unfair treatment of specific groups if the technology is
1833 deployed in applications affecting societal decision-making.

1834
1835

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

A tranquil, winter cabin

A serene, celestial observatory



Figure 24: Visualization of generated objects compared to other edge-cutting methods

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

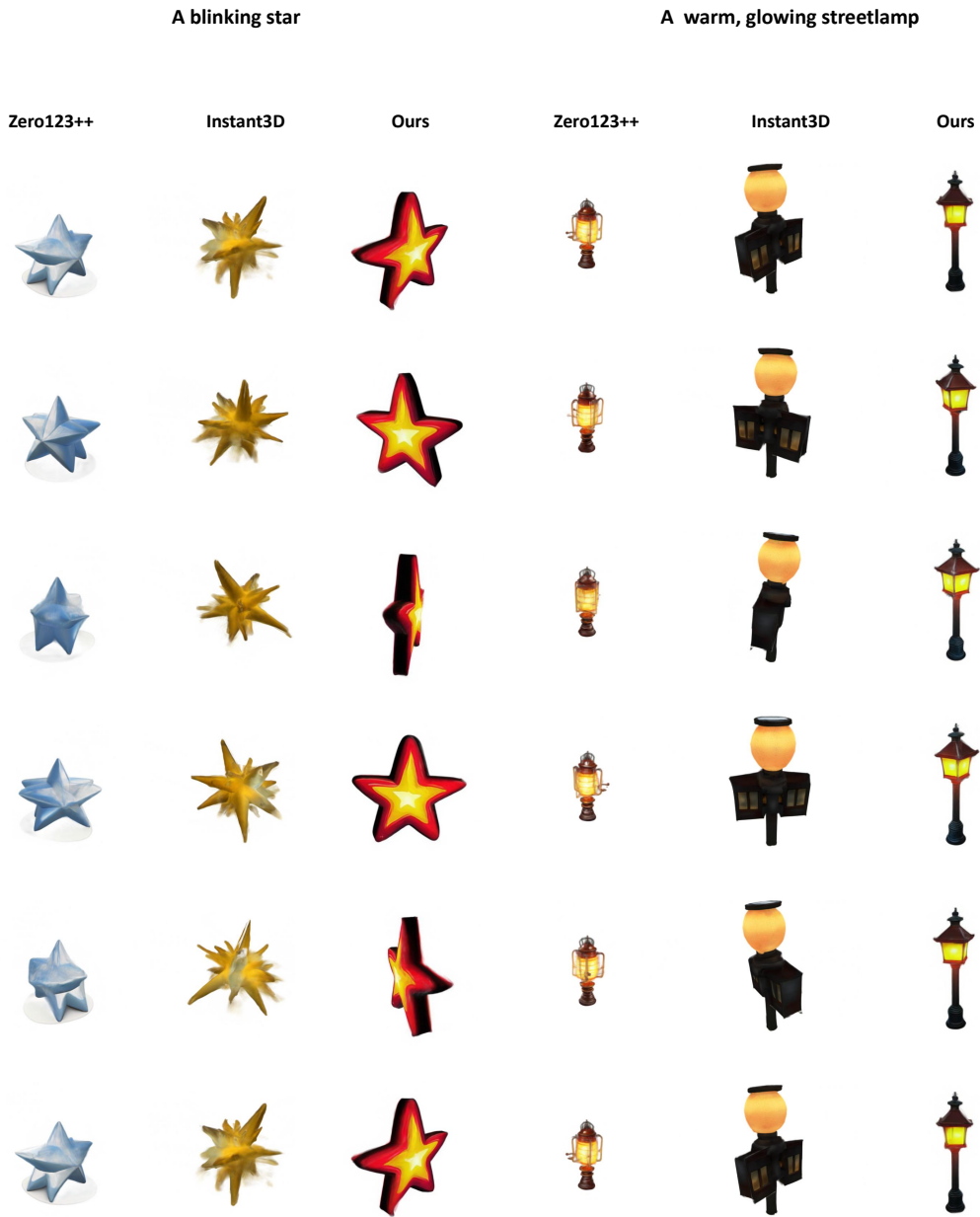


Figure 25: Visualization of generated objects compared to other edge-cutting methods

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

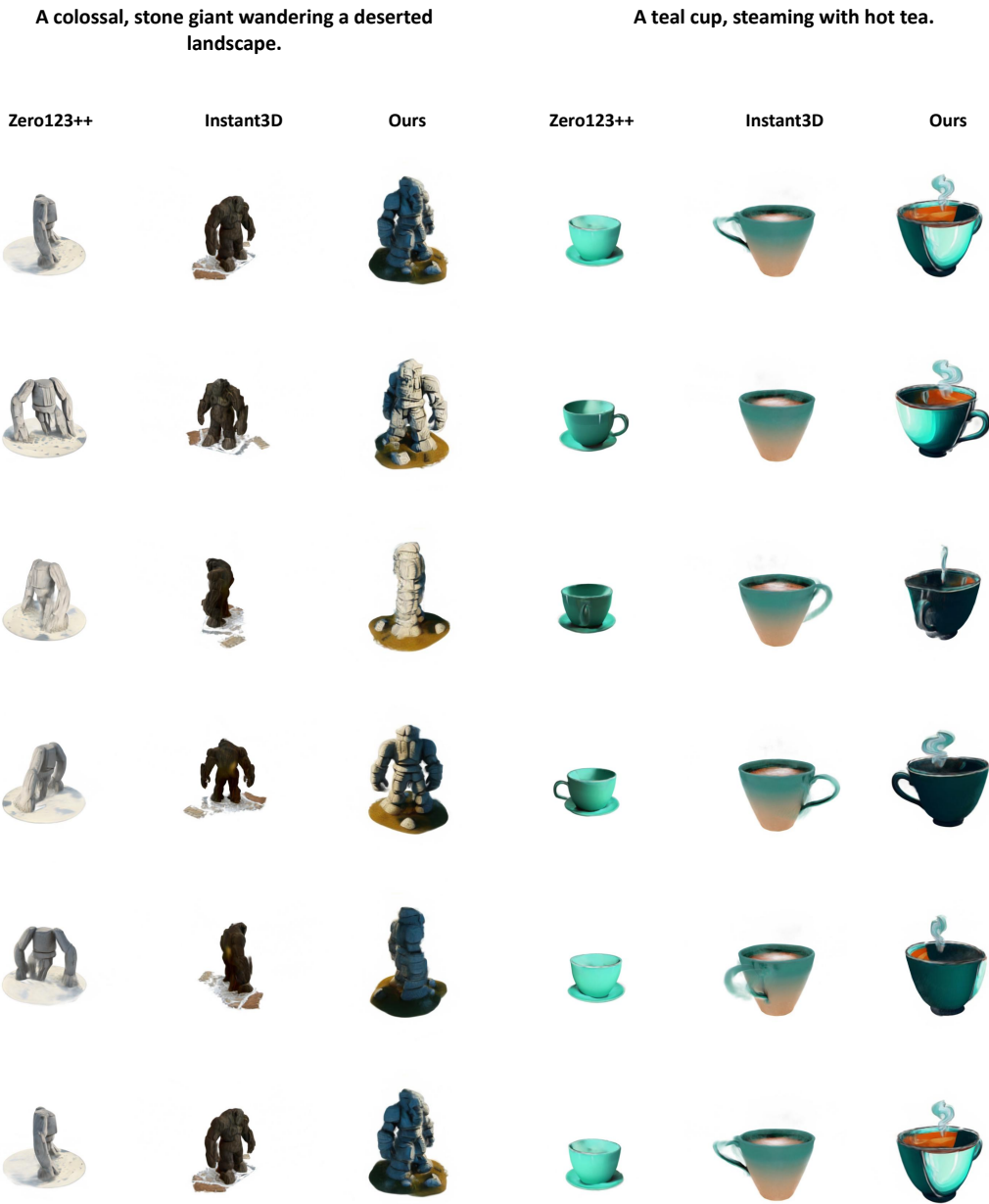


Figure 26: Visualization of generated objects compared to other edge-cutting methods

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051



Figure 27: Visualization of generated objects compared to other edge-cutting methods with different style control.

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105



Figure 28: Visualization of generated objects compared to other edge-cutting methods with different style control.

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159



Figure 29: Visualization of generated objects compared to other edge-cutting methods with different style control.

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213



Figure 30: Visualization of generated objects compared to other edge-cutting methods with different style control.

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267



Figure 31: Visualization of generated objects compared to other edge-cutting methods with different style control.