

# MV-ADAPTER: MULTI-VIEW CONSISTENT IMAGE GENERATION MADE EASY

Anonymous authors

Paper under double-blind review

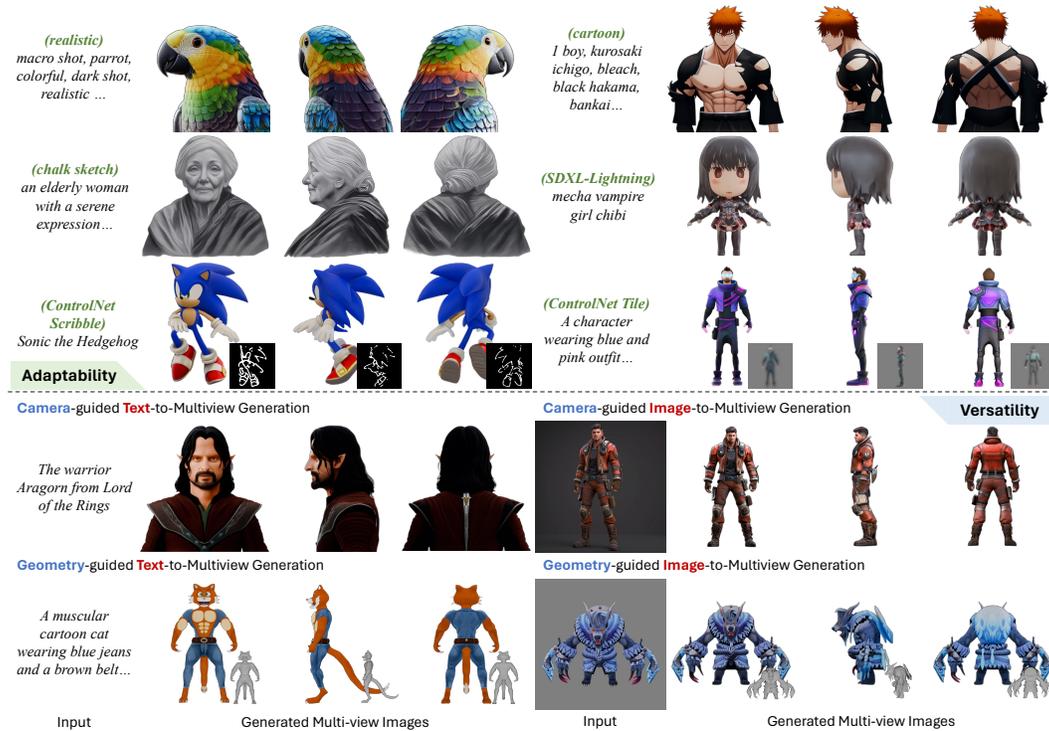


Figure 1: MV-Adapter is a versatile plug-and-play adapter that turns existing pre-trained text-to-image (T2I) diffusion models to multi-view image generators. **Row 1,2,3:** results by integrating MV-Adapter with personalized T2I models, distilled few-step T2I models, and ControlNets (Zhang et al., 2023), demonstrating its **adaptability**. **Row 4,5:** results under various control signals, including view-guided or geometry-guided generation with text or image inputs, showcasing its **versatility**.

## ABSTRACT

Generating multi-view images of an object has important applications in content creation and perception. Existing methods achieved this by making invasive changes to pre-trained text-to-image (T2I) models and performing full-parameter training, leading to three main limitations: (1) High computational costs, especially for high-resolution outputs; (2) Incompatibility with derivatives and extensions of the base model, such as personalized models, distilled few-step models, and plugins like ControlNets; (3) Limited versatility, as they primarily serve a single purpose and cannot handle diverse conditioning signals such as text, images, and geometry. In this paper, we present MV-Adapter, a plug-and-play module working on top of pre-trained T2I models. MV-Adapter enables efficient training for high-resolution synthesis while maintaining full compatibility with all kinds of derivatives of the base T2I model. MV-Adapter provides a unified implementation for generating multi-view images from various conditions, facilitating applications such as text- and image-based 3D generation and texturing. We demonstrate

054 that MV-Adapter sets a new quality standard for multi-view image generation, and  
055 opens up new possibilities due to its adaptability and versatility.  
056

## 057 058 1 INTRODUCTION 059

060 Multi-view image generation is a fundamental task with significant applications in areas such as  
061 2D/3D content creation, robotics perception, and simulation. With the advent of text-to-image (T2I)  
062 diffusion models (Ramesh et al., 2022; Nichol et al., 2022; Saharia et al., 2022; Ramesh et al.,  
063 2021; Balaji et al., 2022; Podell et al., 2024; Mokady et al., 2023), there has been considerable  
064 progress in generating high-quality single-view images. Extending these models to handle multi-  
065 view generation holds the promise of unifying text, image, and 3D data into a cohesive framework.

066 Recent attempts on multi-view image generation (Shi et al., 2023b; Tang et al., 2023; 2024; Huang  
067 et al., 2024b; Gao et al., 2024; Liu et al., 2023a; Long et al., 2024; Li et al., 2024; Kant et al.,  
068 2024; Zheng & Vedaldi, 2024; Wang & Shi, 2023) involve fine-tuning T2I models on large-scale 3D  
069 datasets (Deitke et al., 2023; Yu et al., 2023) and propose modeling 3D consistency across images  
070 by applying self-attention on relevant pixels in different views. As a pioneer work, MVDream (Shi  
071 et al., 2023b) applies self-attention on latent pixels from all generated views, allowing the network  
072 to implicitly learn the consistency. Follow-up works like SPAD (Kant et al., 2024) and Era3D (Li  
073 et al., 2024) constrain the self-attention along epipolar lines, which improves efficiency and enables  
074 higher-resolution synthesis (Li et al., 2024).

075 While these advancements have led to progressively better results, they face several limitations that  
076 hinder their practicality. First, they often require full fine-tuning of pre-trained T2I models, which  
077 demands substantial computational resources and memory usage, making it impractical to scale to  
078 larger models and higher resolutions. The most advanced model to date is trained on Stable Diffusion  
079 2-1 with 860M parameters at resolution 512 (Li et al., 2024). Second, full-parameter training with  
080 substantial network structure changes can lead to catastrophic forgetting of pre-trained knowledge,  
081 impairing compatibility with derivatives and extensions of the base model, including personalized  
082 models tailored to specific subjects or styles (Ruiz et al., 2023; Gal et al., 2022; Hu et al., 2021),  
083 distilled few-step models optimized for efficiency (Luo et al., 2023; Lin et al., 2024), and plugins  
084 (e.g. ControlNets (Zhang et al., 2023)) that add new functionalities. This incompatibility restricts  
085 the ability to leverage the continuous advancements and community contributions. Third, existing  
086 methods mainly serve a single purpose, for example generating multi-view images from text (Shi  
087 et al., 2023b; Kant et al., 2024), a reference image (Wang & Shi, 2023; Shi et al., 2023a; Wang  
088 et al., 2024b; Voleti et al., 2024; Wen et al., 2024; Li et al., 2024; Huang et al., 2024b), or geometry  
089 conditions (Bensadoun et al., 2024), but sharing the underlying logic of maintaining multi-view  
090 consistency. It is desirable to have a unified design that incorporates diverse conditioning signals,  
091 addressing the varied requirements of multi-view generation tasks across various domains.

092 To address these challenges, we propose MV-Adapter, a versatile plug-and-play adapter that en-  
093 hances T2I models and their derivatives for multi-view generation under various conditions. Our  
094 approach eliminates the need for full model fine-tuning by introducing a multi-view adapter network  
095 seamlessly integrated with frozen T2Is. This significantly reduces computational costs and memory  
096 usage in training, making high-resolution generation feasible on larger models like Stable Diffusion  
097 XL (Podell et al., 2024). By preserving the original feature space of the base T2I model during  
098 training, MV-Adapter maintains high compatibility with various derivative models and community-  
099 developed plugins. This adaptability allows users to benefit from personalized subjects or styles, ef-  
100 ficient few-step generation, and additional controllability without specific re-training. Moreover, we  
101 involve a unified design in the adapter network to support diverse conditioning inputs. It comprises  
102 a condition guider that processes camera or geometry guidance, enabling the model to incorporate  
103 viewpoint or structural information and therefore supports both 3D object generation and 3D model  
104 texture generation. This design also introduces decoupled attention blocks, which consists of multi-  
105 view attention layers and optional image cross-attention layers, allowing the model to generate from  
106 both text and image conditions.

107 We evaluate the performance of our MV-Adapter on a diverse set of personalized and efficient T2Is  
from the community. These models encompass a wide spectrum of domains, such as various styles  
and concepts, forming a comprehensive benchmark for our evaluation. Results of our experiments  
demonstrate promising outcomes.

In summary, contributions of MV-Adapter are as follows: (1) **Efficiency**. MV-Adapter eliminates the need for full fine-tuning, increasing training efficiency and enabling high-resolution generation. (2) **Adaptability**. MV-Adapter is fully compatible with derivatives and extensions of the base T2I model. (3) **Versatility**. MV-Adapter supports multiple conditioning inputs, broadening the scope of multi-view generation applications. (4) **Performance**. Experiments demonstrate that T2Is with MV-Adapter can generate multi-view consistent images while preserving visual quality, leveraging the specific strengths of the base T2I models.

## 2 RELATED WORK

**Text-to-image diffusion models.** Text-to-image (T2I) generation (Ramesh et al., 2022; Nichol et al., 2022; Saharia et al., 2022; Ramesh et al., 2021; Balaji et al., 2022; Podell et al., 2024; Mokady et al., 2023; Huang et al., 2024a) has made remarkable progress, particularly with the advancement of diffusion models (Ho et al., 2020; Song et al., 2020; Dhariwal & Nichol, 2021; Ho & Salimans, 2022). Guided diffusion (Dhariwal & Nichol, 2021) and classifier-free guidance (Ho & Salimans, 2022) improved text conditioning and generation fidelity. DALL-E2 (Ramesh et al., 2022) leverages CLIP (Radford et al., 2021) for better text-image alignment. The Latent Diffusion Model (Rombach et al., 2022), also known as Stable Diffusion, enhances efficiency by performing diffusion in the latent space of an autoencoder. Stable Diffusion XL (Podell et al., 2024), a two-stage cascade diffusion model, has greatly improved the generation of high-frequency details and overall image quality, elevating the aesthetic appeal of the outputs.

**Derivatives and extensions of T2I models.** To facilitate creation with pre-trained T2Is, various derivative models and extensions have been developed, focusing on model distillation for efficiency (Meng et al., 2023; Song et al., 2023; Luo et al., 2023; Lin et al., 2024) and controllable generation (Cao et al., 2024). These derivatives encompass personalization (Ruiz et al., 2023; Gal et al., 2022; Hu et al., 2021; Shi et al., 2024; Wang et al., 2024a; Ma et al., 2024; Song et al., 2024; Kumari et al., 2023; Ye et al., 2023), and spatial control (Mou et al., 2024; Zhang et al., 2023). Typically, they employ adapters or fine-tuning methods to extend functionality while preserving the original feature space of the pre-trained models. For instance, DreamBooth (Ruiz et al., 2023) uses class-specific prior preservation loss for personalization, and ControlNet (Zhang et al., 2023) and T2I-Adapter (Mou et al., 2024) enable flexible control over generation by incorporating adapters to the base T2Is. Our work builds on these non-intrusive methods, ensuring compatibility with our MV-Adapter for broader applications.

**Multi-view Generation with T2I models.** Multi-view generation methods (Shi et al., 2023b; Tang et al., 2023; 2024; Huang et al., 2024b; Gao et al., 2024; Liu et al., 2023a; Long et al., 2024; Li et al., 2024; Kant et al., 2024; Zheng & Vedaldi, 2024; Wang & Shi, 2023) extend T2I models by leveraging large-scale 3D datasets (Deitke et al., 2023; Yu et al., 2023). For instance, MVDream (Shi et al., 2023b) integrates camera embeddings and expands the self-attention mechanism from 2D to 3D for cross-view connections, while SPAD (Kant et al., 2024) enhances spatial relational modeling by applying epipolar constraints to cross-view attention. Era3D (Li et al., 2024) introduces an efficient row-wise self-attention mechanism aligned with epipolar lines across views, facilitating high-resolution multi-view generation. However, these methods typically require extensive parameter updates, altering the feature space of pre-trained T2I models and limiting their compatibility with T2I derivatives. Our work addresses this by introducing a multi-view adapter that harmonizes with pre-trained T2Is, significantly expanding the potential for diverse applications.

## 3 PRELIMINARY

Here we introduce the preliminary of multi-view diffusion models (Shi et al., 2023b; Kant et al., 2024; Li et al., 2024), which can help understand the common strategies in modeling multi-view consistency within T2I models.

**Multi-view diffusion models.** Multi-view diffusion models enhance T2Is by introducing multi-view attention mechanism, enabling the generation of images that are consistent across different viewpoints. Several studies (Shi et al., 2023b; Wang & Shi, 2023) extend the self-attention of T2Is to include all pixels across multi-view images. Let  $f^{in}$  denotes the input of the attention block, the dense multi-view self-attention extends  $f^{in}$  from the view itself to the concatenated feature sequence

from  $n$  views. While this approach captures global dependencies, it is computationally intensive, as it processes all pixels of all views. To mitigate the computational cost, epipolar attention (Kant et al., 2024; Huang et al., 2024b) leverages geometric relationships between views. Specifically, methods like SPAD (Kant et al., 2024) extend the self-attention by restricting  $f^{in}$  to the view itself as well as patches along its epipolar lines.

Furthermore, when generating orthographic views at an elevation angle of  $0^\circ$ , the epipolar lines align with the image rows. Utilizing this property, row-wise self-attention (Li et al., 2024) is introduced after the original self-attention layers in T2I models. The process is defined as:

$$f^{self} = \text{SelfAttn}(f^{in}) + f^{in}; f^{mv} = \text{MultiViewAttn}(f^{self}) + f^{self} \quad (1)$$

where MultiViewAttn performs attention across the same rows in different views, effectively enforcing multi-view consistency with reduced computational overhead.

## 4 METHOD

MV-Adapter is a plug-and-play adapter that learns multi-view priors transferable to derivatives of T2Is without specific tuning, and enable them to generate multi-view consistent images under various conditions. As shown in Fig. 2, at inference, our MV-Adapter, which contains a condition guider and the decoupled attention layers, can be inserted into a personalized or distilled T2I to constitute the multi-view generator.

In detail, as shown in Fig. 3, the condition guider in Sec. 4.1 encodes the camera or geometry information, which supports both camera-guided and geometry-guided generation. Within the decoupled attention mechanism in Sec. 4.2, the additional multi-view attention layers learn multi-view consistency, while the optional image cross-attention layers are for image-conditioned generation. Sec. 4.3 elaborates on the training and inference processes of the MV-Adapter.

### 4.1 CONDITION GUIDER

We design a general condition guider that supports encoding both camera and geometric representations, enabling T2I models to perform multi-view generation under various guidance.

**Camera conditioning.** MV-Adapter is designed for generating  $n$  orthographic views. To condition on the camera pose, we use a camera ray representation (“raymap”) that shares the same height and width as the latent representations and encodes the ray origin and direction at each spatial location (Watson et al., 2022; Sajjadi et al., 2022; Gao et al., 2024).

**Geometry conditioning.** Geometry-guided multi-view generation helps applications like texture generation. To condition on the geometry information, we use a global, rather than view-dependent representation that contains position maps and normal maps (Li et al., 2023; Bensadoun et al., 2024). Each pixel in the position map represents the coordinates of the point on the shape, which provide point correspondences across different views. Normal maps provide orientation information and capture fine geometric details, helping produce detailed textures. We concatenate the position map and normal map along to form a composite geometric conditioning input for each view.

**Encoder design.** To encode the camera or geometry representation, we design a simple and lightweight condition guider for the conditioning maps  $c_m$  ( $c_m \in \mathbb{R}^{n \times 6 \times h \times w}$ ). Inspired by T2I-Adapter (Mou et al., 2024), the condition guider consists of a series of convolutional networks, which contain feature extraction blocks and downsampling layers to adapt the feature resolution to the features in the U-Net encoder. The extracted multi-scale features are then added to the corresponding scales in the U-Net, enabling the model to integrate the conditioning information seamlessly at

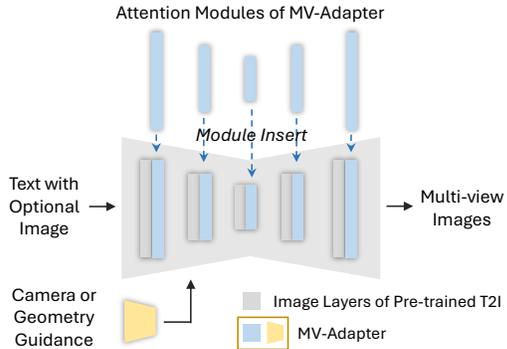


Figure 2: Inference pipeline.

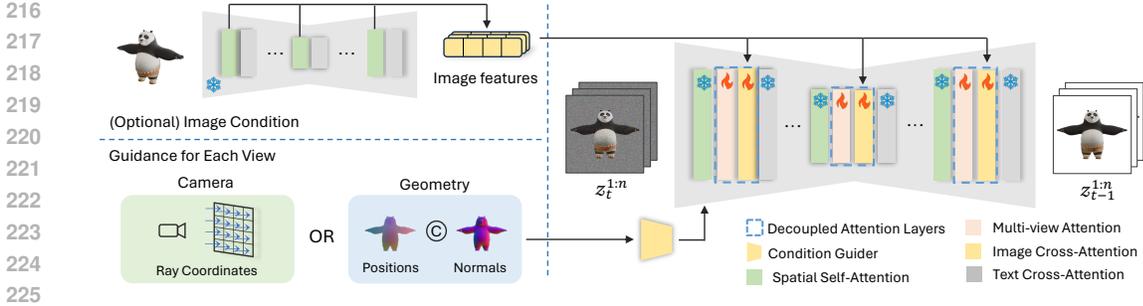


Figure 3: Overview of MV-Adapter. Our MV-Adapter consists of two components: 1) a condition guider that encodes camera or geometry condition; 2) decoupled attention layers that contain multi-view attention for learning multi-view consistency, and optional image cross-attention to support image-conditioned generation, where we use the pre-trained U-Net to encode fine-grained information of the reference image. After training, MV-Adapter can be inserted into any personalized or distilled T2I to generate multi-view images while leveraging the specific strengths of base models.

multiple levels. In theory, the input to our encoder is not limited to specific types of conditions; it can also be extended to a wider variety of maps, such as depth maps and pose maps.

#### 4.2 DECOUPLED ATTENTION

We introduce a decoupled attention mechanism, where we retain the original spatial self-attention layers and add multi-view attention layers that enforce multi-view consistency as well as optional image cross-attention layers for image-conditioned generation. These three types of attention layers are organized in a parallel architecture, effectively leveraging the image priors from the pre-trained self-attention layers.

**Multi-view attention.** Considering the different applications of camera-guided and geometry-guided multi-view generation, we design different strategies for multi-view attention to meet the specific needs of each application (shown in Fig. 4(a)). For camera-guided generation, we follow Era3D (Li et al., 2024) to achieve image-to-3D creation, allowing the model to generate multi-view images at an elevation of  $0^\circ$ . We then employ row-wise self-attention, restricting the multi-view attention to process only patches within the same row across views. For geometry-guided generation, considering the view coverage requirements of its main application (*i.e.*, texture generation), we adjust the distribution of the generated multi-view images. In addition to the four views evenly at elevation  $0^\circ$ , we add two views from top and bottom. We perform both row-wise and column-wise self-attention, enabling efficient information exchange among all views.

**Image cross-attention.** To condition on reference images  $c_i$  and achieve control over fine-grained appearance details, we propose a novel method for incorporating detailed information from the image without altering the original feature space of the T2I model. We employ the pre-trained T2I model itself as the image encoder. Specifically, we employ a frozen U-Net that is identical to the pre-trained SD U-Net (Rombach et al., 2022), with its weights initialized from the SD U-Net. During the feature extraction process, we pass the clear reference image into this frozen U-Net, setting the timestep  $t = 0$ , and then extract multi-scale features from the spatial self-attention layers. These fine-grained features contain detailed information about the subject and are injected into the denoising U-Net through the decoupled image cross-attention layers. In this way, we leverage the rich representations learned by the pre-trained model, enabling precise control over the generated content.

**Attention architecture.** In the pre-trained T2I model, the spatial self-attention layer and text cross-attention layer are connected serially through residual connections. Suppose feature sequence  $f^{in}$  is the input of the attention block, we can express the process as

$$f^{self} = \text{SelfAttn}(f^{in}) + f^{in}; f^{cross} = \text{CrossAttn}(f^{self}) + f^{self} \quad (2)$$

A straightforward method to incorporate new attention layers is to append them after the original layers, connecting them in a serial manner. However, the sequential arrangement may not effectively

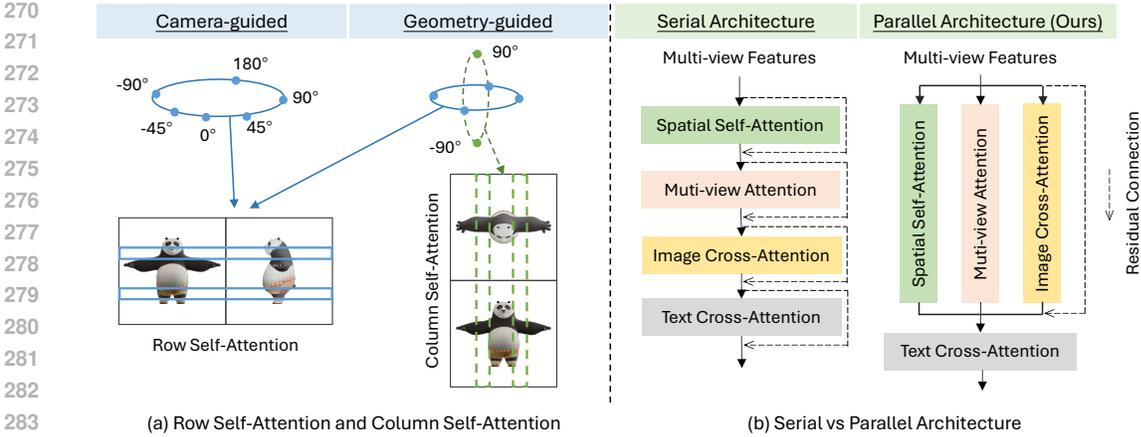


Figure 4: Overview of the decoupled attention design. (a) For camera-guided generation, similar to Era3D (Li et al., 2024), we apply row-wise self-attention to generate multi-view images at an elevations of  $0^\circ$ . For geometry-guided generation, designed for texture generation, we add two views from the top and bottom to ensure comprehensive coverage and perform both row-wise and column-wise self-attention. (b) Instead of serially connecting new attention layers, which requires training additional modules from scratch, we utilize a parallel architecture that builds upon the established priors of pre-trained self-attention, enabling more efficient learning.

utilize the image priors modeled by the pre-trained self-attention layers, as it requires the new layers to learn from scratch. To fully exploit the effective priors of the spatial self-attention layers, we adopt a parallel architecture, as shown in Fig. 4(b). The process can be formulated as

$$\mathbf{f}^{self} = \text{SelfAttn}(\mathbf{f}^{in}) + \text{MultiViewAttn}(\mathbf{f}^{in}) + \text{ImageCrossAttn}(\mathbf{f}^{in}, \mathbf{f}^{ref}) + \mathbf{f}^{in} \quad (3)$$

where  $\mathbf{f}^{ref}$  refers to features of the reference image. Since the features  $\mathbf{f}^{in}$  fed into the new layers are the same as those to the self-attention layer, we can effectively initialize them with the pre-trained layers to transfer the image priors. We zero-initialize the output projection layer of the new layers to ensure that the initial output does not disrupt the original feature space. This architectural choice allows the model to build upon the established priors, facilitating efficient learning of multi-view consistency and image-conditioned generation, while preserving the original space of the base T2Is.

### 4.3 TRAINING AND INFERENCE

During training, we only optimize the MV-Adapter, while freezing weights of the pre-trained T2I models. We train MV-Adapter on the dataset with pairs of a reference image, text and  $n$  views, using the same training objective as T2I models:

$$\mathcal{L} = \mathbb{E}_{\mathcal{E}(\mathbf{x}_0^{1:n}), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{c}_t, \mathbf{c}_i, \mathbf{c}_m, t} [\|\epsilon - \epsilon_\theta(\mathbf{z}_t^{1:n}, \mathbf{c}_t, \mathbf{c}_i, \mathbf{c}_m, t)\|_2^2] \quad (4)$$

where  $\mathbf{c}_t$ ,  $\mathbf{c}_i$  and  $\mathbf{c}_m$  represent texts, reference images and conditioning maps (i.e., camera or geometry conditions) respectively. We randomly zero out the features of the reference image to drop image conditions, enabling classifier-free guidance at inference. Similar to prior work (Blattmann et al., 2023; Hoogetboom et al., 2023), we shift the noise schedule towards high noise levels as we move from the T2Is to the multi-view diffusion model that captures data of higher dimensionality. We shift the log signal-to-noise ratio by  $\log(n)$ , where  $n$  is the number of generated views.

## 5 EXPERIMENTS

We implemented MV-Adapter on Stable Diffusion V2.1 (SD2.1) and Stable Diffusion XL (SDXL), training a  $512 \times 512$  adapter for SD2.1 and a  $768 \times 768$  adapter for SDXL using a subset of the Objaverse dataset (Deitke et al., 2023). Detailed configurations are provided in the Appendix.



Figure 5: Results with community models and extensions. Each sample corresponds to a distinct T2I model or extension. Information about the models can be found in the Appendix.

Table 1: Quantitative comparison on camera-guided text-to-multiview generation.

Method	FID↓	IS↑	CLIP Score↑
MVDream	32.15	14.38	31.76
SPAD	48.79	12.04	30.87
Ours (SD2.1)	31.24	15.01	32.04
Ours (SDXL)	<b>29.71</b>	<b>16.38</b>	<b>33.17</b>

Table 2: Quantitative comparison on camera-guided image-to-multiview generation.

Method	PSNR↑	SSIM↑	LPIPS↓
ImageDream	19.280	0.8472	0.1218
Zero123++	20.312	0.8417	0.1205
CRM	20.185	0.8325	0.1247
SV3D	20.042	0.8267	0.1396
Ouroboros3D	20.810	0.8535	0.1193
Era3D	20.890	0.8601	0.1199
Ours (SD2.1)	20.867	0.8695	0.1147
Ours (SDXL)	<b>22.131</b>	<b>0.8816</b>	<b>0.1002</b>

### 5.1 CAMERA-GUIDED MULTI-VIEW GENERATION

**Evaluation on community models and extensions.** We evaluated MV-Adapter using representative T2Is and extensions, including personalized models (Ruiz et al., 2023; Hu et al., 2021), efficient distilled models (Luo et al., 2023; Lin et al., 2024), and plugins such as ControlNet (Zhang et al., 2023). We present eight qualitative results in Fig. 5. More results can be found in the Appendix.

**Comparison with baselines.** For text-to-multiview generation, we compared our MV-Adapter with MVDream (Shi et al., 2023b) and SPAD (Kant et al., 2024) on 1,000 prompts from the Objaverse dataset. The results are presented in Fig. 6 and Table 1. For image-to-multiview generation, we conduct comparison with ImageDream (Wang & Shi, 2023), Zero123++ (Shi et al., 2023a), CRM (Wang et al., 2024b), SV3D (Voleti et al., 2024), Ouroboros3D (Wen et al., 2024), and Era3D (Li et al., 2024) on the Google Scanned Objects (GSO) dataset (Downs et al., 2022), as results shown in Fig. 7 and Table 2. Experiments indicate that, by preserving the original feature space of T2I models, our MV-Adapter achieves higher visual fidelity and consistency with conditions.



Figure 6: Qualitative comparison on camera-guided text-to-multiview generation. our MV-Adapter achieves higher visual fidelity and image-text consistency.



Figure 7: Qualitative comparison on camera-guided image-to-multiview generation.

## 5.2 GEOMETRY-GUIDED MULTI-VIEW GENERATION

**Evaluation on community models and extensions.** We evaluated our geometry-guided model with T2I derivative models. The results in Fig. 8 demonstrate the adaptability of MV-Adapter in seamlessly integrating with different base models.

**Comparison with baselines.** We compare our text- and image-conditioned multi-view-based texture generation method (see Sec. 5.4) with four state-of-the-art methods, including TEXTure (Richardson et al., 2023), Text2Tex (Chen et al., 2023), Paint3D (Zeng et al., 2024), SyncMVD (Liu et al., 2023b), and FlashTex (Deng et al., 2024). For our image-to-texture model, we used ControlNet (Zhang et al., 2023) to generate reference images conditioned on text and depth maps. As shown in Fig. 10 and Table 3, compared to these project-and-inpaint or synchronized multi-view texturing methods, our approach fine-tunes additional modules to model geometric associations and preserves the generative capabilities of the base T2I model, thereby producing multi-view consistent and high-quality textures. Additionally, testing on a single RTX 4090 GPU revealed that our method achieves faster generation speeds than the others.

Table 3: Quantitative comparison on 3D texture generation. FID and KID ( $\times 10^{-4}$ ) are evaluated on multi-view renderings. Our models achieves best texture quality with faster inference.

Method	FID↓	KID↓	Time↓
TEXTure	56.44	61.16	90s
Text2Tex	58.43	60.81	421s
Paint3D	44.38	47.06	60s
SyncMVD	36.13	42.28	50s
FlashTex	50.48	56.36	186s
Ours (SD2.1 - Text)	38.19	42.83	<b>18s</b>
Ours (SD2.1 - Image)	33.93	38.73	19s
Ours (SDXL - Text)	32.75	35.18	32s
Ours (SDXL - Image)	<b>27.28</b>	<b>29.47</b>	33s

## 5.3 ABLATION STUDY

We conduct ablation studies to evaluate the efficiency and adaptability of our MV-Adapter, as well as the detailed design of the adapter network.



449 Figure 8: Results of geometry-guided text-to-multiview generation with community models.

452 **Efficiency.** To assess the training efficiency of our adapter design, we conducted comparison with Era3D (Li et al., 2024), which requires full training rather than fine-tuning only adapters like us. We further extend this model to SDXL (Podell et al., 2024) for a comprehensive evaluation. As shown in Table 4, our MV-Adapter significantly reduces training costs, facilitating high-resolution multi-view generation based on larger backbones.

452 Table 4: Comparison of training costs with full-tuning methods (batch size set to 1).

Method	Trainable params ↓	Memory usage ↓	Training speed ↑
Era3D (SD2.1)	993M	36G	2.2iter/s
Ours (SD2.1)	<b>127M</b>	<b>17G</b>	<b>3.1iter/s</b>
Era3D (SDXL)	3.1B	>80G	-
Ours (SDXL)	<b>490M</b>	<b>60G</b>	<b>1.05iter/s</b>

463 **Adaptability.** We compare MV-Adapter with the full-trained text-to-multiview generation method MVDream (Shi et al., 2023b) regarding compatibility with T2I derivatives. MVDream, which fine-tunes the whole T2I model, cannot be easily replaced with other T2Is; thus, we integrate LoRA (Hu et al., 2021) for our experiments. As shown in Fig. 9, MVDream struggles to generate images that align with the text and style, whereas our MV-Adapter produces high-quality results, demonstrating its superior adaptability.

464 *(3d style) 1 girl, blue eyes, upper body, mask, eyes half closed*

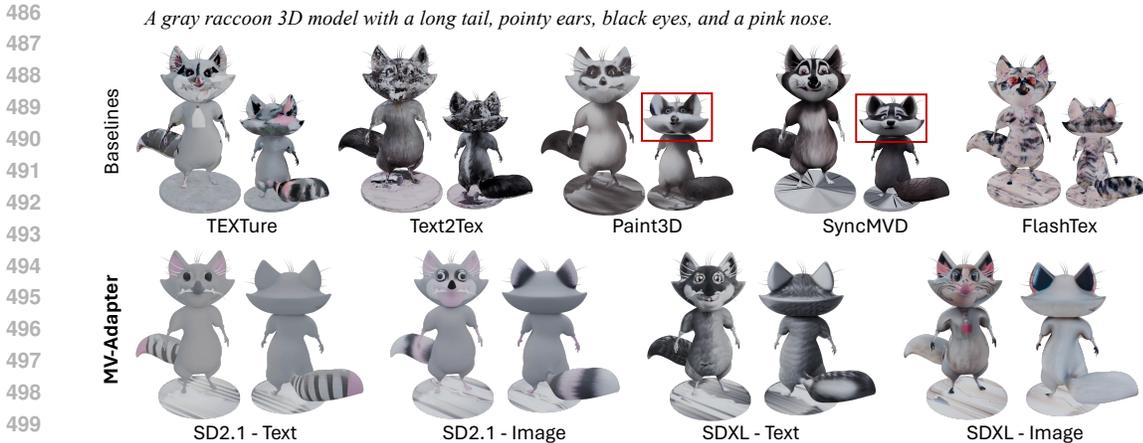


474 Figure 9: Qualitative ablation study on the adaptability of MV-Adapter.

475 **Network design.** We conducted ablation studies on our proposed image encoder and parallel attention architecture. Specifically, we compare the settings of a) using CLIP (Radford et al., 2021) for encoding reference images instead of SD U-Net, and b) replacing the parallel architecture with a serial counterpart, with c) our MV-Adapter. As shown in Fig. 11, comparing a) and c) reveals that CLIP capture only coarse, semantic-level information, while the pre-trained U-Net encodes finer details, producing results closely aligned with the input. Comparing b) and c) shows that, the serial setting, which does not leverage the pre-trained image prior, tends to produce artifacts and misaligned details. Our MV-Adapter achieves greater consistency both among generated views and with the reference image, especially at the detail level. More results can be found in the Appendix.

476 Table 5: Quantitative ablation studies on attention architecture.

Method	PSNR ↑	SSIM ↑	LPIPS ↓
Serial (SDXL)	20.687	0.8681	0.1149
Parallel (SDXL)	<b>22.131</b>	<b>0.8816</b>	<b>0.1002</b>



501 Figure 10: Qualitative comparison on texture generation. We compare our text- and image-  
502 conditioned models with baseline methods.



510 Figure 11: Qualitative ablation study on the network design.

511

512

513 5.4 APPLICATIONS

514

515 **3D generation.** We follow the existing pipelines (Li et al., 2024; Wu et al., 2024) to achieve 3D  
516 generation. After generating multi-view images from text or image conditions using MV-Adapter,  
517 we use StableNormal (Ye et al., 2024) to generate corresponding normal maps. The multi-view  
518 images and normal maps are then fed into NeuS (Wang et al., 2021) to reconstruct the 3D mesh.  
519 The generated results are shown in the Appendix.

520

521 **Texture generation.** We use backprojection and incidence-based weighted blending tech-  
522 niques (Bensadoun et al., 2024) to map the generated multi-view images onto the UV texture map.  
523 Despite optimizing view distribution to enhance coverage, some areas may remain uncovered due  
524 to occlusions or extreme angles. To address this, we perform view coverage analysis to identify  
525 uncovered regions, render images from the current 3D texture for those views, and refine them using  
526 an efficient inpainting model (Suvorov et al., 2022). We show more visual results in the Appendix.

527

528 6 CONCLUSION

529

530 In this paper, we introduce MV-Adapter, a versatile plug-and-play adapter that enhances text-to-  
531 image (T2I) diffusion models and their derivatives for multi-view generation under various condi-  
532 tions, without compromising quality or modifying the original feature space. Our approach incor-  
533 porates a condition guider and a decoupled attention mechanism, enabling both camera-guided and  
534 geometry-guided multi-view generation from text and images. Once trained, our MV-Adapter can  
535 be seamlessly integrated into various T2I models—including personalized, distilled, and plugin-  
536 enhanced models—to generate multi-view images with high consistency and visual fidelity. Exten-  
537 sive evaluations highlight the efficiency, adaptability, and versatility of MV-Adapter across different  
538 models and conditions. Furthermore, we extend our multi-view generation framework to support  
539 applications such as 3D generation and texture generation. Overall, MV-Adapter offers an efficient  
and flexible solution for multi-view image generation, significantly broadening the capabilities of  
pre-trained T2I models and presenting exciting possibilities for a wide range of applications.

## REFERENCES

- 540  
541  
542 Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten  
543 Kreis, Miika Aittala, Timo Aila, Samuli Laine, et al. ediff-i: Text-to-image diffusion models with  
544 an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.
- 545 Raphael Bensadoun, Yanir Kleiman, Idan Azuri, Omri Harosh, Andrea Vedaldi, Natalia Neverova,  
546 and Oran Gafni. Meta 3d texturegen: Fast and consistent texture generation for 3d objects. *arXiv*  
547 *preprint arXiv:2407.02430*, 2024.
- 548 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendeleevitch, Maciej Kilian, Dominik  
549 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling  
550 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- 551  
552 Pu Cao, Feng Zhou, Qing Song, and Lu Yang. Controllable generation with text-to-image diffusion  
553 models: A survey. *arXiv preprint arXiv:2403.04279*, 2024.
- 554  
555 Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner.  
556 Text2tex: Text-driven texture synthesis via diffusion models. In *Proceedings of the IEEE/CVF*  
557 *International Conference on Computer Vision*, pp. 18558–18568, 2023.
- 558  
559 Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig  
560 Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of anno-  
561 tated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
562 *Recognition*, pp. 13142–13153, 2023.
- 563  
564 Kangle Deng, Timothy Omernick, Alexander Weiss, Deva Ramanan, Jun-Yan Zhu, Tinghui Zhou,  
565 and Maneesh Agrawala. Flashtex: Fast relightable mesh texturing with lightcontrolnet. *arXiv*  
566 *preprint arXiv:2402.13251*, 2024.
- 567  
568 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*  
569 *in neural information processing systems*, 34:8780–8794, 2021.
- 570  
571 Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann,  
572 Thomas B McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset  
573 of 3d scanned household items. In *2022 International Conference on Robotics and Automation*  
574 *(ICRA)*, pp. 2553–2560. IEEE, 2022.
- 575  
576 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel  
577 Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual  
578 inversion. *arXiv preprint arXiv:2208.01618*, 2022.
- 579  
580 Ruiqi Gao, Aleksander Holynski, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla, Pratul  
581 Srinivasan, Jonathan T Barron, and Ben Poole. Cat3d: Create anything in 3d with multi-view  
582 diffusion models. *arXiv preprint arXiv:2405.10314*, 2024.
- 583  
584 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*  
585 *arXiv:2207.12598*, 2022.
- 586  
587 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
588 *neural information processing systems*, 33:6840–6851, 2020.
- 589  
590 Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for  
591 high resolution images. In *International Conference on Machine Learning*, pp. 13213–13232.  
592 PMLR, 2023.
- 593  
594 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,  
595 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*  
596 *arXiv:2106.09685*, 2021.
- 597  
598 Zehuan Huang, Hongxing Fan, Lipeng Wang, and Lu Sheng. From parts to whole: A unified  
599 reference framework for controllable human image generation. *arXiv preprint arXiv:2404.15267*,  
600 2024a.

- 594 Zehuan Huang, Hao Wen, Junting Dong, Yaohui Wang, Yangguang Li, Xinyuan Chen, Yan-Pei  
595 Cao, Ding Liang, Yu Qiao, Bo Dai, et al. Epidiff: Enhancing multi-view synthesis via localized  
596 epipolar-constrained diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
597 and Pattern Recognition*, pp. 9784–9794, 2024b.
- 598  
599 Yash Kant, Aliaksandr Siarohin, Ziyi Wu, Michael Vasilkovsky, Guocheng Qian, Jian Ren, Riza Alp  
600 Guler, Bernard Ghanem, Sergey Tulyakov, and Igor Gilitschenski. Spad: Spatially aware multi-  
601 view diffusers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
602 Recognition*, pp. 10026–10038, 2024.
- 603 Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept  
604 customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Com-  
605 puter Vision and Pattern Recognition*, pp. 1931–1941, 2023.
- 606  
607 Black Forest Labs. Flux. [Online], 2024. [https://github.com/black-forest-labs/  
608 flux](https://github.com/black-forest-labs/flux).
- 609 Peng Li, Yuan Liu, Xiaoxiao Long, Feihu Zhang, Cheng Lin, Mengfei Li, Xingqun Qi, Shanghang  
610 Zhang, Wenhan Luo, Ping Tan, et al. Era3d: High-resolution multiview diffusion using efficient  
611 row-wise attention. *arXiv preprint arXiv:2405.11616*, 2024.
- 612  
613 Weiyu Li, Rui Chen, Xuelin Chen, and Ping Tan. Sweetdreamer: Aligning geometric priors in 2d  
614 diffusion for consistent text-to-3d. *arXiv preprint arXiv:2310.02596*, 2023.
- 615  
616 Shanchuan Lin, Anran Wang, and Xiao Yang. Sd-xl-lightning: Progressive adversarial diffusion  
617 distillation. *arXiv preprint arXiv:2402.13929*, 2024.
- 618  
619 Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang.  
620 Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint  
621 arXiv:2309.03453*, 2023a.
- 622  
623 Yuxin Liu, Minshan Xie, Hanyuan Liu, and Tien-Tsin Wong. Text-guided texturing by synchronized  
624 multi-view diffusion. *arXiv preprint arXiv:2311.12891*, 2023b.
- 625  
626 Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma,  
627 Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. Wonder3d: Single image to 3d  
628 using cross-domain diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
629 and Pattern Recognition*, pp. 9970–9980, 2024.
- 630  
631 Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthe-  
632 sizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
- 633  
634 Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with pre-  
635 trained models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 636  
637 Jian Ma, Junhao Liang, Chen Chen, and Haonan Lu. Subject-diffusion: Open domain personalized  
638 text-to-image generation without test-time fine-tuning. In *ACM SIGGRAPH 2024 Conference  
639 Papers*, pp. 1–12, 2024.
- 640  
641 Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and  
642 Tim Salimans. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF  
643 Conference on Computer Vision and Pattern Recognition*, pp. 14297–14306, 2023.
- 644  
645 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for  
646 editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF Conference  
647 on Computer Vision and Pattern Recognition*, pp. 6038–6047, 2023.
- 648  
649 Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan.  
650 T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffu-  
651 sion models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(5):4296–4304,  
652 Mar. 2024. doi: 10.1609/aaai.v38i5.28226. URL [https://ojs.aaai.org/index.php/  
653 AAAI/article/view/28226](https://ojs.aaai.org/index.php/AAAI/article/view/28226).

- 648 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob  
649 McGrew, Ilya Sutskever, and Mark Chen. GLIDE: towards photorealistic image generation and  
650 editing with text-guided diffusion models. In *International Conference on Machine Learning,  
651 ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine  
652 Learning Research*, pp. 16784–16804, 2022.
- 653 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of  
654 the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- 655  
656 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
657 Penna, and Robin Rombach. SDXL: improving latent diffusion models for high-resolution image  
658 synthesis. In *The Twelfth International Conference on Learning Representations, ICLR 2024,  
659 Vienna, Austria, May 7-11, 2024*, 2024.
- 660 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
661 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
662 models from natural language supervision. In *International conference on machine learning*, pp.  
663 8748–8763. PMLR, 2021.
- 664  
665 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,  
666 and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine  
667 learning*, pp. 8821–8831. Pmlr, 2021.
- 668  
669 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-  
670 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- 671  
672 Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-  
673 guided texturing of 3d shapes. In *ACM SIGGRAPH 2023 conference proceedings*, pp. 1–11,  
674 2023.
- 675  
676 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
677 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
678 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 679  
680 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-  
681 ical image segmentation. In *Medical image computing and computer-assisted intervention—  
682 MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceed-  
683 ings, part III 18*, pp. 234–241. Springer, 2015.
- 684  
685 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
686 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Pro-  
687 ceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22500–  
688 22510, 2023.
- 689  
690 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
691 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic  
692 text-to-image diffusion models with deep language understanding. *Advances in neural informa-  
693 tion processing systems*, 35:36479–36494, 2022.
- 694  
695 Mehdi SM Sajjadi, Henning Meyer, Etienne Pot, Urs Bergmann, Klaus Greff, Noha Radwan, Suhani  
696 Vora, Mario Lučić, Daniel Duckworth, Alexey Dosovitskiy, et al. Scene representation trans-  
697 former: Geometry-free novel view synthesis through set-latent scene representations. In *Proceed-  
698 ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6229–6238,  
699 2022.
- 700  
701 Viraj Shah, Nataniel Ruiz, Forrester Cole, Erika Lu, Svetlana Lazebnik, Yuanzhen Li, and Varun  
Jampani. Ziplora: Any subject in any style by effectively merging loras. *arXiv preprint  
arXiv:2311.13600*, 2023.
- Jing Shi, Wei Xiong, Zhe Lin, and Hyun Joon Jung. Instantbooth: Personalized text-to-image gen-  
eration without test-time finetuning. In *Proceedings of the IEEE/CVF Conference on Computer  
Vision and Pattern Recognition*, pp. 8543–8552, 2024.

- 702 Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen,  
703 Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base  
704 model. *arXiv preprint arXiv:2310.15110*, 2023a.
- 705
- 706 Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view  
707 diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023b.
- 708
- 709 Yukai Shi, Jianan Wang, He Cao, Boshi Tang, Xianbiao Qi, Tianyu Yang, Yukun Huang, Shilong  
710 Liu, Lei Zhang, and Heung-Yeung Shum. Toss: High-quality text-guided novel view synthesis  
711 from a single image. *arXiv preprint arXiv:2310.10644*, 2023c.
- 712
- 713 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv  
preprint arXiv:2010.02502*, 2020.
- 714
- 715 Kunpeng Song, Yizhe Zhu, Bingchen Liu, Qing Yan, Ahmed Elgammal, and Xiao Yang. Moma:  
716 Multimodal llm adapter for fast personalized image generation. *arXiv preprint arXiv:2404.05674*,  
717 2024.
- 718
- 719 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint  
arXiv:2303.01469*, 2023.
- 720
- 721 Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha,  
722 Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky.  
723 Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the  
724 IEEE/CVF winter conference on applications of computer vision*, pp. 2149–2159, 2022.
- 725
- 726 Shitao Tang, Fuyang Zhang, Jiacheng Chen, Peng Wang, and Yasutaka Furukawa. Mvdifffusion:  
727 Enabling holistic multi-view image generation with correspondence-aware diffusion. *arXiv*, 2023.
- 728
- 729 Shitao Tang, Jiacheng Chen, Dilin Wang, Chengzhou Tang, Fuyang Zhang, Yuchen Fan, Vikas  
730 Chandra, Yasutaka Furukawa, and Rakesh Ranjan. Mvdifffusion++: A dense high-resolution  
731 multi-view diffusion model for single or sparse-view 3d object reconstruction. *arXiv preprint  
arXiv:2402.12712*, 2024.
- 732
- 733 Kolors Team. Kolors: Effective training of diffusion model for photorealistic text-to-image synthe-  
734 sis. *arXiv preprint*, 2024.
- 735
- 736 Vikram Voleti, Chun-Han Yao, Mark Boss, Adam Letts, David Pankratz, Dmitry Tochilkin, Chris-  
737 tian Laforte, Robin Rombach, and Varun Jampani. Sv3d: Novel multi-view synthesis and 3d  
738 generation from a single image using latent video diffusion. *arXiv preprint arXiv:2403.12008*,  
2024.
- 739
- 740 Haofan Wang, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. Instantstyle: Free lunch towards  
741 style-preserving in text-to-image generation. *arXiv preprint arXiv:2404.02733*, 2024a.
- 742
- 743 Peng Wang and Yichun Shi. Imagedream: Image-prompt multi-view diffusion for 3d generation.  
*arXiv preprint arXiv:2312.02201*, 2023.
- 744
- 745 Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus:  
746 Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv  
preprint arXiv:2106.10689*, 2021.
- 747
- 748 Zhengyi Wang, Yikai Wang, Yifei Chen, Chendong Xiang, Shuo Chen, Dajiang Yu, Chongxuan Li,  
749 Hang Su, and Jun Zhu. Crm: Single image to 3d textured mesh with convolutional reconstruction  
750 model. *arXiv preprint arXiv:2403.05034*, 2024b.
- 751
- 752 Daniel Watson, William Chan, Ricardo Martin-Brualla, Jonathan Ho, Andrea Tagliasacchi, and Mo-  
753 hammad Norouzi. Novel view synthesis with diffusion models. *arXiv preprint arXiv:2210.04628*,  
2022.
- 754
- 755 Hao Wen, Zehuan Huang, Yaohui Wang, Xinyuan Chen, Yu Qiao, and Lu Sheng. Ouroboros3d:  
Image-to-3d generation via 3d-aware recursive diffusion. *arXiv preprint arXiv:2406.03184*, 2024.

- 756 Kailu Wu, Fangfu Liu, Zhihan Cai, Runjie Yan, Hanyang Wang, Yating Hu, Yueqi Duan, and  
 757 Kaisheng Ma. Unique3d: High-quality and efficient 3d mesh generation from a single image.  
 758 *arXiv preprint arXiv:2405.20343*, 2024.
- 759  
 760 Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,  
 761 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models  
 762 with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.
- 763  
 764 Chongjie Ye, Lingteng Qiu, Xiaodong Gu, Qi Zuo, Yushuang Wu, Zilong Dong, Liefeng Bo, Yuliang  
 765 Xiu, and Xiaoguang Han. Stablenormal: Reducing diffusion variance for stable and sharp normal.  
*ACM Transactions on Graphics (TOG)*, 2024.
- 766  
 767 Hu Ye, Jun Zhang, Sibao Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt  
 768 adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.
- 769  
 770 Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan,  
 771 Chenming Zhu, Zhangyang Xiong, Tianyou Liang, et al. Mvimngnet: A large-scale dataset of  
 772 multi-view images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern  
 773 recognition*, pp. 9150–9161, 2023.
- 774  
 775 Xianfang Zeng, Xin Chen, Zhongqi Qi, Wen Liu, Zibo Zhao, Zhibin Wang, Bin Fu, Yong Liu, and  
 776 Gang Yu. Paint3d: Paint anything 3d with lighting-less texture diffusion models. In *Proceedings  
 777 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4252–4262, 2024.
- 778  
 779 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
 pp. 3836–3847, 2023.
- 780  
 781 Chuanxia Zheng and Andrea Vedaldi. Free3d: Consistent novel view synthesis without 3d represen-  
 782 tation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
 pp. 9720–9731, 2024.
- 783  
 784 Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun  
 785 Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all,  
 786 March 2024. URL <https://github.com/hpcaitech/Open-Sora>.

## 788 A APPENDIX

### 789 A.1 BACKGROUND

790  
 791 **Stable Diffusion (SD) and Stable Diffusion XL (SDXL).** We adopt Stable Diffusion (Rombach  
 792 et al., 2022) and Stable Diffusion XL (Podell et al., 2024) as our base T2I models, since they have a  
 793 well-developed community with many powerful derivatives for evaluation. SD and SDXL perform  
 794 the diffusion process within the latent space of a pre-trained autoencoder  $\mathcal{E}(\cdot)$  and  $\mathcal{D}(\cdot)$ . In training,  
 795 an encoded image  $z_0 = \mathcal{E}(x_0)$  is perturbed to  $z_t$  at step  $t$  by the forward diffusion. The denoising  
 796 network  $\epsilon_\theta$  learns to reverse this process by predicting the added noise, encouraged by an MSE loss:  
 797

$$798 \mathcal{L} = \mathbb{E}_{\mathcal{E}(x_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), c, t} [\|\epsilon - \epsilon_\theta(z_t, c, t)\|_2^2] \quad (5)$$

799 where  $c$  denotes the conditioning texts. In SD,  $\epsilon_\theta$  is implemented as a UNet (Ronneberger et al.,  
 800 2015) consisting of pairs of down/up sample blocks and a middle block. Each block contains pairs  
 801 of spatial self-attention layers and cross-attention layers, which are serially connected using the  
 802 residual structure. SDXL leverages a three times larger UNet backbone than SD for high-resolution  
 803 image synthesis, and introduces a refinement denoiser to improve the visual fidelity.  
 804

### 805 A.2 IMPLEMENTATION DETAILS

806  
 807 **Dataset.** We trained MV-Adapter on a filtered high-quality subset of the Objaverse dataset (Deitke  
 808 et al., 2023), comprising approximately 70,000 samples, with captions from Cap3D (Luo et al.,  
 809 2024). To accommodate the efficient multi-view self-attention mechanism, we rendered ortho-  
 graphic views to train the the model to generate  $n = 6$  views per sample. For the camera-guided

810 generation, we rendered views of 3D models with the elevation angle set to  $0^\circ$  and azimuth angles  
 811 at  $\{0^\circ, 45^\circ, 90^\circ, 180^\circ, 270^\circ, 315^\circ\}$ . This distribution aligns with the setting used in Era3D (Li et al.,  
 812 2024), facilitating the application of a similar image-to-3D pipeline for 3D generation tasks. For  
 813 the geometry-guided generation, we included four views at an elevation of  $0^\circ$  with azimuth angles  
 814 of  $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$ , added two additional views from the top and bottom. In addition to the  
 815 target views, we rendered five random views within a certain frontal range of the models to serve as  
 816 reference images during training.

817 **Training.** We utilized two versions of Stable Diffusion (Rombach et al., 2022) as the base mod-  
 818 els for training. Specifically, we trained a 512-resolution model based on Stable Diffusion 2.1  
 819 (SD2.1) and a 768-resolution model based on Stable Diffusion XL (SDXL). During training, we  
 820 randomly dropped the text condition with a probability of 0.1, the image condition with a probabili-  
 821 ty of 0.1, and both text and image conditions simultaneously with a probability of 0.1. Following  
 822 prior work (Hooeboom et al., 2023; Blattmann et al., 2023), we shifted the noise schedule to higher  
 823 noise levels by adjusting the log signal-to-noise ratio (SNR) by  $\log(n)$ , where  $n = 6$  is the number  
 824 of the generated views. For the specific training configurations, we used a learning rate of  $5 \times 10^{-5}$   
 825 and trained the MV-Adapter on 8 NVIDIA A100 GPUs for 10 epochs.

827 **Inference.** In our experimental setup, we generated multi-view images using the DDPM sam-  
 828 pler (Ho et al., 2020) with classifier-free guidance (Ho & Salimans, 2022), and set the number of  
 829 inference steps to 50. For generation conditioned solely on text (i.e., setting the weight of the image  
 830 condition  $\lambda_i$  to 0), we set the guidance scale to 7.0. For image-conditioned generation, we set the  
 831 guidance scale of image condition  $\alpha$  and text condition  $\beta$  to 3.0. Following TOSS (Shi et al., 2023c),  
 832 the calculation can be expressed as:

$$833 \hat{\epsilon}_\theta(\mathbf{z}_t^{1:n}, \mathbf{c}_t, \mathbf{c}_i, \mathbf{c}_m, t) = \epsilon_\theta(\mathbf{z}_t^{1:n}, \emptyset, \emptyset, \mathbf{c}_m, t) \\
 834 + \alpha [\epsilon_\theta(\mathbf{z}_t^{1:n}, \emptyset, \mathbf{c}_i, \mathbf{c}_m, t) - \epsilon_\theta(\mathbf{z}_t^{1:n}, \emptyset, \emptyset, \mathbf{c}_m, t)] \\
 835 + \beta [\epsilon_\theta(\mathbf{z}_t^{1:n}, \mathbf{c}_t, \mathbf{c}_i, \mathbf{c}_m, t) - \epsilon_\theta(\mathbf{z}_t^{1:n}, \emptyset, \mathbf{c}_i, \mathbf{c}_m, t)] \quad (6)$$

836 where  $\mathbf{c}_t$ ,  $\mathbf{c}_i$  and  $\mathbf{c}_m$  represent texts, reference images and conditioning maps (i.e., camera or geom-  
 837 etry conditions) respectively. Since we did not drop  $\mathbf{c}_m$  during the training process, we do not use  
 838 the classifier-free guidance method for it.

840 **Comparison with baselines.** We conducted comprehensive comparisons with baseline methods  
 841 across three settings: text-to-multiview generation, image-to-multiview generation, and texture gen-  
 842 eration. In these experiments, we evaluated both versions of MV-Adapter based on Stable Diffusion  
 843 2.1 (SD2.1) (Rombach et al., 2022) and Stable Diffusion XL (SDXL) (Podell et al., 2024), demon-  
 844 strating the performance gains brought by MV-Adapter due to its efficient training and scalability.

845 For text-to-multiview generation, we selected MVDream (Shi et al., 2023b) and SPAD (Kant et al.,  
 846 2024) as baseline methods. MVDream extends the original self-attention mechanism of T2I models  
 847 to the multi-view domain. SPAD introduces epipolar constraints into the multi-view attention mech-  
 848 anism. We tested on 1,000 prompts selected from the Objaverse dataset (Deitke et al., 2023). We  
 849 computed Fréchet Inception Distance (FID), Inception Score (IS), and CLIP Score on all generated  
 850 views to assess the quality of the generated images and their alignment with the textual prompts.

851 For image-to-multiview generation, we compared our method with ImageDream (Wang & Shi,  
 852 2023), Zero123++(Shi et al., 2023a), CRM(Wang et al., 2024b), SV3D (Voleti et al., 2024),  
 853 Ouroboros3D (Wen et al., 2024), and Era3D (Li et al., 2024). ImageDream, Zero123++, CRM,  
 854 and Era3D generally fall into the category of modifying the original network architecture of T2I  
 855 models to extend them for multi-view generation. SV3D and Ouroboros3D fine-tune text-to-video  
 856 (T2V) models to achieve multi-view generation. We selected 100 assets covering multiple object  
 857 categories from the Google Scanned Objects (GSO) dataset (Downs et al., 2022) as our test set. For  
 858 each asset, we rendered input images from front-facing views, with input views randomly distributed  
 859 in azimuth angles between  $-45^\circ$  and  $45^\circ$  and elevation angles between  $-10^\circ$  and  $30^\circ$ . We evalu-  
 860 ated the generated multi-view images by computing Peak Signal-to-Noise Ratio (PSNR), Structural  
 861 Similarity Index Measure (SSIM), and Learned Perceptual Image Patch Similarity (LPIPS) between  
 862 the generated images and the ground truth, assessing both the consistency and quality of the outputs.

863 For 3D texture generation, we compared our text-based and image-based models with project-  
 and-paint methods such as TEXTure (Richardson et al., 2023), Text2Tex (Chen et al., 2023), and

Table 6: Community models and extensions for evaluation.

Category	Model Name	Domain	Model Type
Personalized T2I	Dreamshaper <sup>1</sup>	General	T2I Base Model
	RealVisXL <sup>2</sup>	Realistic	T2I Base Model
	Animagine-xl <sup>3</sup>	2D Cartoon	T2I Base Model
	3D Render Style XL <sup>4</sup>	3D Cartoon	LoRA
	Pokemon Trainer Sprite PixelArt <sup>5</sup>	Pixel Art	LoRA
	Chalk Sketch SDXL <sup>6</sup>	Chalk Sketch	LoRA
	Chinese Ink LoRA <sup>7</sup>	Color Ink	LoRA
	Zen Ink Wash Sumi-e <sup>8</sup>	Wash Ink	LoRA
	Watercolor Style SDXL <sup>9</sup>	Watercolor	LoRA
	Papercut SDXL <sup>10</sup>	Papercut	LoRA
	Furry Enhancer <sup>11</sup>	Enhancer	LoRA
	White Pitbull Dog SDXL <sup>12</sup>	Concept	LoRA
	Spider spirit fourth sister <sup>13</sup>	Concept	LoRA
Distilled T2I	SDXL-Lightning <sup>14</sup>	Few Step	T2I Base Model
	LCM-SDXL <sup>15</sup>	Few Step	T2I Base Model
Extension	ControlNet Openpose <sup>16</sup>	Spatial Control	Plugin
	ControlNet Scribble <sup>17</sup>	Spatial Control	Plugin
	ControlNet Tile <sup>18</sup>	Image Deblur	Plugin
	T2I-Adapter Sketch <sup>19</sup>	Spatial Control	Plugin
	IP-Adapter <sup>20</sup>	Image Prompt	Plugin

Paint3D (Zeng et al., 2024), the synchronized multi-view texturing method SyncMVD (Liu et al., 2023b), and the optimization-based method FlashTex (Deng et al., 2024). We randomly selected 200 models along with their captions from the Objaverse (Deitke et al., 2023) dataset for testing. Multiple views were rendered from the generated 3D textures, and we computed FID and Kernel Inception Distance (KID) of them to evaluate the quality of the generated textures. Additionally, we recorded the texture generation time to assess the inference efficiency of each method.

**Community models and extensions for evaluation.** To ensure a comprehensive benchmark, we selected a diverse set of representative T2I derivative models and extensions from the community

<sup>1</sup><https://civitai.com/models/112902?modelVersionId=126688>

<sup>2</sup><https://civitai.com/models/139562?modelVersionId=789646>

<sup>3</sup><https://huggingface.co/cagliostrolab/animagine-xl-3.1>

<sup>4</sup>[https://huggingface.co/goofyai/3d\\_render\\_style\\_xl](https://huggingface.co/goofyai/3d_render_style_xl)

<sup>5</sup><https://civitai.com/models/159333/pokemon-trainer-sprite-pixelart?modelVersionId=443092>

<sup>6</sup><https://huggingface.co/JerryOrbachJr/Chalk-Sketch-SDXL>

<sup>7</sup>[https://huggingface.co/ming-yang/sdxl\\_chinese\\_ink\\_lora](https://huggingface.co/ming-yang/sdxl_chinese_ink_lora)

<sup>8</sup><https://civitai.com/models/647926/zen-ink-wash-sumi-e-sdxl-pony-flux?modelVersionId=724876>

<sup>9</sup><https://civitai.com/models/484723/watercolor-style-sdxl>

<sup>10</sup>[https://huggingface.co/TheLastBen/Papercut\\_SDXL](https://huggingface.co/TheLastBen/Papercut_SDXL)

<sup>11</sup><https://civitai.com/models/310964/furry-enhancer?modelVersionId=558568>

<sup>12</sup><https://civitai.com/models/700883/white-pitbull-dog-sdxl?modelVersionId=787948>

<sup>13</sup><https://civitai.com/models/689010/pony-black-myth-wukong-spider-spirit-fourth-sister?modelVersionId=771146>

<sup>14</sup><https://huggingface.co/ByteDance/SDXL-Lightning>

<sup>15</sup><https://huggingface.co/latent-consistency/lcm-sdxl>

<sup>16</sup><https://huggingface.co/xinsir/controlnet-openpose-sdxl-1.0>

<sup>17</sup><https://huggingface.co/xinsir/controlnet-scribble-sdxl-1.0>

<sup>18</sup><https://huggingface.co/xinsir/controlnet-tile-sdxl-1.0>

<sup>19</sup><https://huggingface.co/TencentARC/t2i-adapter-sketch-sdxl-1.0>

<sup>20</sup><https://huggingface.co/h94/IP-Adapter>

918 for evaluation. As illustrated in Table 6, these models include personalized models that encompass  
 919 various domains such as anime, stylistic paintings, and realistic photographic images, as well as ef-  
 920 ficient distilled models and plugins for controllable generation. They cover a wide range of subjects,  
 921 including portraits, animals, landscapes, and more. This selection enables a thorough evaluation of  
 922 our approach across different styles and content, demonstrating the adaptability and generality of  
 923 MV-Adapter in working with various T2I derivatives and extensions.

### 924 925 926 A.3 ADDITIONAL DISCUSSIONS

#### 927 928 A.3.1 MV-ADAPTER VS. MULTI-VIEW LoRA

929  
930 LoRA (Low-Rank Adaptation) (Hu et al., 2021) offers an alternative approach to achieving plug-  
 931 and-play multi-view generation. Specifically, using a condition encoder to inject camera represen-  
 932 tations, we extend the original self-attention mechanism to operate across all pixels of multiple  
 933 views. During training, we introduce trainable LoRA layers into the network, allowing these lay-  
 934 ers to learn multi-view consistency or, optionally, generate images conditioned on a reference view.  
 935 This approach requires the spatial self-attention mechanism to simultaneously capture spatial image  
 936 knowledge, ensure multi-view consistency, and align generated images with reference views.

937 However, the multi-view LoRA approach has a notable limitation. The “incremental changes” it  
 938 introduces to the network are **not orthogonal or decoupled** from those induced by T2I derivatives,  
 939 such as personalized T2I models or LoRAs. Specifically, layers fine-tuned by multi-view LoRA  
 940 and those tuned by personalized LoRA often overlap. Note that each weight matrix learned by both  
 941 represents a linear transformation defined by its columns, so it is intuitive that the merger would  
 942 retain the information available in these columns only when the columns that are being added are  
 943 orthogonal to each other (Shah et al., 2023). Clearly, the multi-view LoRA and personalized models  
 944 are not orthogonal, which often leads to challenges in retaining both sets of learned knowledge. This  
 945 can result in a trade-off where either multi-view consistency or the fidelity of concepts (such as style  
 946 or subject identity) is compromised.

947 In contrast, our proposed **decoupled** attention mechanism encourages different attention layers to  
 948 specialize in their respective tasks without needing to fine-tune the original spatial self-attention  
 949 layers. In this design, the layers we train do not overlap with those in the original T2I model,  
 950 thereby better preserving the original feature space and enhancing compatibility with other models.

951 We conducted a series of experiments to test these approaches. We trained two versions of multi-  
 952 view LoRA, targeting different modules: (1) inserting LoRA layers only into the attention layers, and  
 953 (2) inserting LoRA layers into multiple layers, including the convolutional layers, down-sampling,  
 954 up-sampling layers, etc. For both settings, we set the LoRA rank to 64 and alpha to 32. As shown in  
 955 Fig. 12 and Fig. 13, while the multi-view LoRA approach can generate multi-view consistent images  
 956 when the base model is not changed, it often struggles to maintain multi-view consistency when  
 957 switching to a different base model or when integrating a new LoRA. In contrast, as demonstrated in  
 958 Fig. 14, our MV-Adapter, equipped with the decoupled attention mechanism, maintains consistent  
 959 multi-view generation even when used with personalized models.

960 Compared to the LoRA mechanism, our decoupled attention-based approach proves more robust  
 961 and adaptable for extending T2I models to multi-view generation, offering greater flexibility and  
 962 compatibility with various pre-trained models.

#### 963 964 A.3.2 ADAPTABILITY OF IMAGE-CONDITIONED MODEL

965  
966 Evaluating the adaptability of the image-conditioned MV-Adapter on personalized models poses a  
 967 challenge because the reference image already provides detailed subject-specific appearance guid-  
 968 ance for multi-view generation. As a result, it’s difficult to assess how well the model adapts when  
 969 the subject’s details are pre-defined. To address this, we conducted experiments on efficient distilled  
 970 models, such as SDXL-Lighting (Lin et al., 2024). As illustrated in Fig. 15, after replacing the base  
 971 model with a distilled T2I variant, the MV-Adapter was able to generate high-quality and multi-view  
 consistent images **in just four steps**.



993 Figure 12: Results of multi-view LoRA (set target modules to attention layers). The azimuth angles  
 994 of the images from left to right are  $0^\circ$ ,  $45^\circ$ ,  $90^\circ$ ,  $180^\circ$ ,  $270^\circ$ ,  $315^\circ$ , corresponding to the front,  
 995 front-left, left, back, right, and front-right of the object.  
 996  
 997



1018 Figure 13: Results of multi-view LoRA (set target modules to attention layers, convolutional layers,  
 1019 etc.). The azimuth angles of the images from left to right are  $0^\circ$ ,  $45^\circ$ ,  $90^\circ$ ,  $180^\circ$ ,  $270^\circ$ ,  $315^\circ$ , corre-  
 1020 sponding to the front, front-left, left, back, right, and front-right of the object.  
 1021  
 1022  
 1023

1024 The experiments clearly demonstrate that our image-conditioned MV-Adapter exhibits strong adapt-  
 1025 ability. Even when integrated into distilled models, it is capable of rapidly generating high-quality  
 multi-view images, proving its efficiency and versatility.



1047 Figure 14: Results of MV-Adapter, which introduces decoupled attention mechanism rather than  
 1048 LoRA. The azimuth angles of the images from left to right are  $0^\circ$ ,  $45^\circ$ ,  $90^\circ$ ,  $180^\circ$ ,  $270^\circ$ ,  $315^\circ$ , cor-  
 1049 responding to the front, front-left, left, back, right, and front-right of the object.



1062 Figure 15: Results of MV-Adapter on camera-guided image-to-multiview generation with SDXL-  
 1063 Lightning (Lin et al., 2024) (number of inference steps set to 4).

### 1066 A.3.3 IMAGE RESTORATION CAPABILITIES

1067

1068 During the training of MV-Adapter, we probabilistically compress the resolution of reference im-  
 1069 ages in the training data pairs to enhance the robustness of multi-view generation from images. We  
 1070 observed that the model trained with this approach is capable of generating high-resolution, de-  
 1071 tailed multi-view images even when the input is low-resolution, as depicted in Fig. 16. Through  
 1072 such training strategy, MV-Adapter has inherent image restoration capabilities and automatically  
 1073 enhances and refines input images during the generation process.

### 1075 A.3.4 SERIAL VS. PARALLEL ATTENTION ARCHITECTURE

1076

1077 To assess the effectiveness of our proposed parallel attention architecture, we conducted ablation  
 1078 studies on image-to-multi-view generation setting. As shown in Fig. 17, the serial setting, which  
 1079 cannot leverage the pre-trained image prior, tends to produce artifacts and inconsistent details with  
 the image input. In contrast, our parallel setting produces high-quality and highly consistent results.



**Inspiration for related tasks.** Our MV-Adapter represents a successful practice of decoupling image priors from geometric knowledge within T2I diffusion models. This approach provides valuable insights for downstream tasks that rely on image priors but also require modeling of geometric, physical, or temporal aspects. Specifically, characteristics related to geometry and viewpoint—such as zooming in/out, lighting variations, and shadow dynamics—can potentially be addressed by introducing new layers that decouple these factors or by fine-tuning the multi-view attention layers of MV-Adapter. By extending this decoupled architecture, it may be possible to model geometric-related properties more effectively, enabling advancements in areas like view-dependent appearance synthesis, relighting, and even animation where temporal consistency is crucial. This opens avenues for future research to explore how similar strategies can be applied to disentangle and control other complex factors in image generation tasks.

### A.3.6 EXTENDING MV-ADAPTER FOR ARBITRARY VIEW SYNTHESIS

In the main text, we introduced a novel adapter architecture—comprising parallel attention layers and a unified condition encoder—to achieve multi-view generation. We implemented efficient row-wise and column-wise attention mechanisms tailored for two specific applications: 3D object generation and 3D texture mapping, generating six views accordingly. However, our adapter framework is not limited to these configurations and can be extended to perform arbitrary view synthesis. To explore this capability, we designed a corresponding approach and conducted experiments, training a new version of MV-Adapter to handle arbitrary viewpoints.

Following CAT3D (Gao et al., 2024), we perform multiple rounds of multi-view generation, with the number of views generated each time set to  $n = 8$ . Starting from text or an initial single image as input, we first generate eight anchor views that broadly cover the object. In practice, these anchor views are positioned at elevations of  $0^\circ$  and  $30^\circ$ , with azimuth angles evenly distributed around the circle (e.g. every  $45^\circ$ ). For generating new target views, we cluster the viewpoints based on their spatial orientations, grouping them into clusters of 8. We then select the 4 nearest known views from the already generated anchor views to serve as conditions guiding the generation of each target view.

In terms of implementation, the overall framework of our MV-Adapter remains unchanged. We adjust its inputs and specific attention components to accommodate arbitrary view synthesis. First, we set the number of input images to either 1 or 4. When using four input views, we concatenate them into a long image and input this into the pre-trained T2I U-Net to extract features. This simple yet effective method allows the images from the four views to interact within the pre-trained U-Net without requiring additional camera embeddings to represent these views. Second, we utilize full self-attention in the multi-view attention component, expanding the attention scope to enable the generation of target views with more flexible distributions.

To train an MV-Adapter capable of generating arbitrary viewpoints, we rendered data from 40 different views, with elevations of  $-10^\circ, 0^\circ, 10^\circ, 20^\circ, 30^\circ$ , and azimuth angles evenly distributed around 360 degrees at each elevation layer. We trained the model for 16 epochs. During the first 8 epochs, the model was trained using a setting of one conditional view and eight target anchor views. In the subsequent 8 epochs, we trained with an equal mixture of one condition plus eight target views and four conditions plus eight target views.

As shown in Fig. 18, the visualization results demonstrate that MV-Adapter can generate consistent, high-quality multi-view images beyond the six views designed for specific applications. This extension further verifies the scalability and practicality of our adapter framework, showcasing its potential for arbitrary view synthesis in diverse applications. More results can be found in the supplementary materials.

### A.4 LIMITATIONS AND FUTURE WORKS

**Domain gap between synthetic data and natural images.** A domain gap exists between the synthetic multi-view data rendered from 3D datasets (Deitke et al., 2023) and natural images, particularly in terms of background presence and visual fidelity. The model trained with synthetic data will be affected to some extent by the specific 3D style appearance, which may affect the generalization of the model. Although the adapter design successfully leverages the priors from the pre-trained T2I model, the quality of the generated images is still influenced by the suboptimal visual quality of the training data. A potential solution involves augmenting the training data with real video datasets,

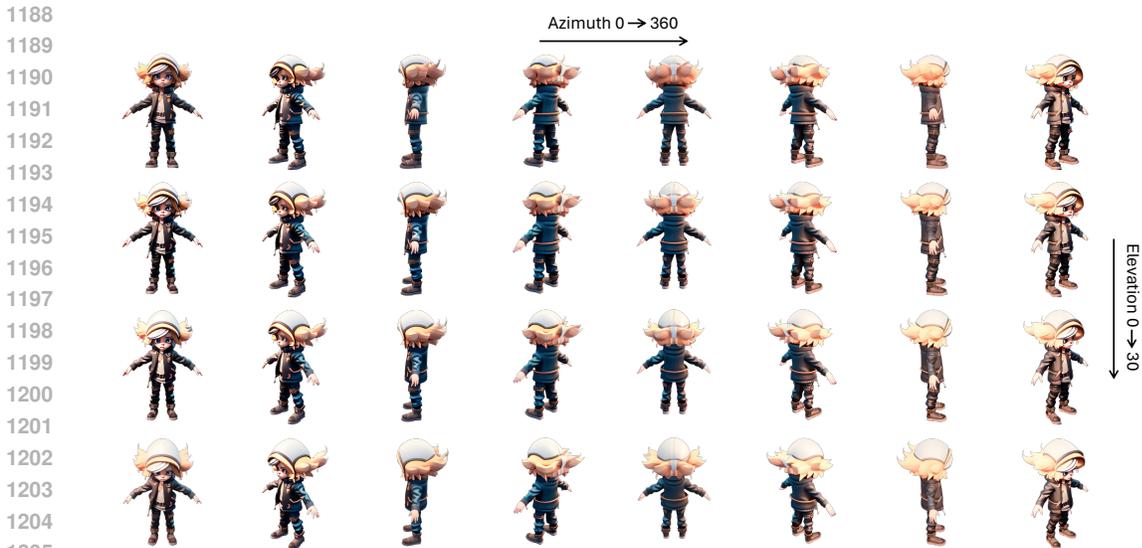


Figure 18: Visualization results using MV-Adapter to generate arbitrary viewpoints.

such as MVImgNet (Yu et al., 2023), which could reduce the domain gap. Additionally, during inference, we recommend incorporating a reference image as additional content control, which will improve the visual fidelity the controllability of the multi-view generation.

**Dependency on image backbone.** Within our decoupled attention mechanism, the visual content, multi-view consistency and alignment with the reference image originate from the underlying image backbone, multi-view attention, and image cross-attention mechanisms, respectively. Notably, both the multi-view attention and image cross-attention layers are initialized using the parameters of the original spatial self-attention layers. Consequently, the overall performance of MV-Adapter is heavily dependent on the base T2I model. If the foundational model struggles to generate content that aligns with the provided prompt or produces images of low quality, MV-Adapter is unlikely to compensate for these deficiencies. On the other hand, employing superior image backbones can enhance the synthetic results. We present a comparison of outputs generated using SDXL (Podell et al., 2024) and SD2.1 (Rombach et al., 2022) models in Fig. 19, which confirms this observation, particularly in text-conditioned multi-view generation. We believe that MV-Adapter can be further developed by utilizing advanced T2I models (Team, 2024; Labs, 2024) based on the DiT architecture (Peebles & Xie, 2023), to achieve higher visual quality in the generated images.



Figure 19: Qualitative comparison of our MV-Adapter based on SD2.1 and SDXL.

**Future works: sparse-view input, 3D scene generation, dynamic multi-view video generation.** This paper provides extensive analyses and enhancements for our novel multi-view adapter, MV-

Adapter. While our model has significantly improved efficiency, adaptability, versatility, and performance compared to previous models, we identify three promising areas for future work:

- Sparse-view input. To enhance controllability in multi-view generation, we can input sparse views into our image encoder (*i.e.*, pre-trained SD U-Net), allowing multiple views to guide the generation process.
- 3D scene generation. We conducted experiments on synthetic data. Our method can be extended to scene-level multi-view generation, accommodating both camera- and geometry-guided approaches with text or image conditions.
- Dynamic multi-view video generation. Exploring dynamic multi-view video generation using a similar approach as MV-Adapter within text-to-video generation models (Zheng et al., 2024; Yang et al., 2024) presents a valuable opportunity for further advancements.

**Future works: modeling new knowledge like MV-Adapter.** By decoupling the learning of geometric knowledge from the image prior, our framework efficiently integrates new knowledge without compromising the base model’s rich visual capabilities. This principle enhances learning from limited data and inspires other tasks that build upon existing image priors to learn new types of knowledge. Beyond multi-view consistency, our approach can be extended to learn zoom in/out effects, consistent lighting conditions, and other viewpoint-dependent attributes. It is possible to model viewpoint-dependent attributes such as lighting, shadows, and reflections by fine-tuning our decoupled multi-view attention on some specific small datasets, which can be defined as personalization or customization of geometric knowledge. MV-Adapter also provides insights for modeling physical or temporal knowledge based on image priors, paving the way for future research in related domains.

## A.5 MORE COMPARISON RESULTS

### A.5.1 IMAGE-TO-MULTI-VIEW GENERATION

To provide a more in-depth analysis of our quantitative results on image-to-multi-view generation, we conducted a user study comparing MV-Adapter (based on SD2.1 (Rombach et al., 2022)) with baseline methods (Wang & Shi, 2023; Shi et al., 2023a; Wang et al., 2024b; Voleti et al., 2024; Wen et al., 2024; Li et al., 2024). The study aimed to evaluate both multi-view consistency and image quality preferences. We selected 30 samples covering a diverse range of categories, such as toy cars, medicine bottles, stationery, dolls, and sculptures. A total of 50 participants were recruited to provide their preferences between the outputs of different methods.

Participants were presented with pairs of multi-view images generated by MV-Adapter and the baseline methods. For each pair, they were asked to choose the one they preferred in terms of multi-view consistency and image quality. The results of the user study are summarized in Fig. 20. The findings indicate that, in terms of multi-view consistency, MV-Adapter performs comparably to Era3D, with preference rates of 25.07% and 22.33%, respectively. However, regarding image quality, MV-Adapter demonstrates a significant advantage, receiving a higher preference rate of 36.80% compared to the baseline methods. The improved image quality can be attributed to MV-Adapter’s ability to leverage the strengths of the underlying T2I models without full fine-tuning, preserving the original feature space and benefiting from the high-quality priors of the base models.

Additionally, we provide supplementary qualitative comparison results in Fig. 21, showcasing side-by-side examples of images generated by MV-Adapter and the baseline methods. These examples further illustrate the superior image quality and consistency achieved by MV-Adapter, highlighting finer details, better texture reproduction, and more coherent structures across different views.

### A.5.2 IMAGE-TO-3D GENERATION

To further evaluate the consistency of multi-view generation and the applicability of MV-Adapter to downstream tasks, we conducted a quantitative comparison of 3D reconstruction performance using MV-Adapter and Era3D (Li et al., 2024), which shares a similar pipeline with our method. The comparison was performed on the Google Scanned Objects (GSO) dataset, focusing on metrics such as Chamfer Distance and Volumetric IoU to assess the geometric quality of the reconstructed 3D models.

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

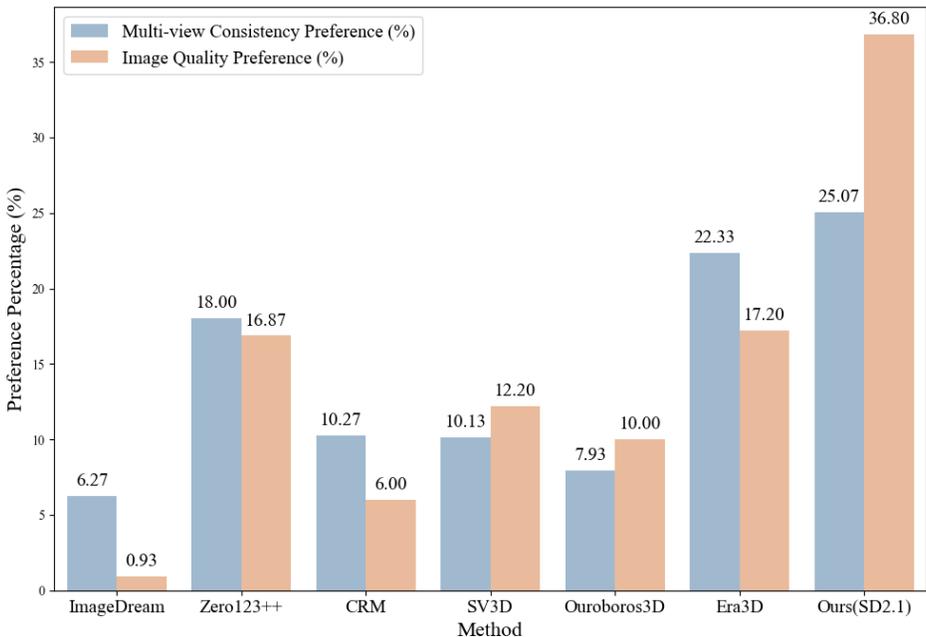


Figure 20: Results of user study on image-to-multi-view generation.

The results, summarized in Table 7, show that 3D reconstruction quality using MV-Adapter based on Stable Diffusion 2.1 (SD2.1) is comparable to that achieved with Era3D. However, when using MV-Adapter based on Stable Diffusion XL (SDXL), the reconstruction quality is significantly higher, with notable improvements in both Chamfer Distance and Volumetric IoU. This demonstrates that MV-Adapter’s efficient training design facilitates compatibility with larger and more advanced base models, such as SDXL, thereby delivering superior results in 3D reconstruction tasks. These findings underline the scalability of MV-Adapter and its ability to leverage the strengths of state-of-the-art T2I models, providing additional benefits to downstream tasks like 3D generation.

Table 7: Quantitative comparison on 3D reconstruction.

Method	Chamfer Distance↓	Volume IoU↑
Era3D	0.0329	0.5118
Ours (SD2.1)	0.0317	0.5173
Ours (SDXL)	<b>0.0206</b>	<b>0.5682</b>

### A.6 MORE VISUAL RESULTS

In Fig. 22 and Fig. 23, we show more visual results of MV-Adapter on camera-guided text-to-multiview generation with community models and extensions, such as ControlNet (Zhang et al., 2023) and IP-Adapter (Ye et al., 2023). In Fig. 24, we show more visual results on camera-guided image-to-multiview generation. In Fig. 25, we show more visual results on text-to-3D generation. In Fig. 26, we show more visual results on image-to-3D generation. In Fig. 27, we show more visual results on geometry-guided text-to-texture generation. In Fig. 28, we show more visual results on geometry-guided image-to-texture generation. Note that we have removed the background of the generated images in the visual results.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

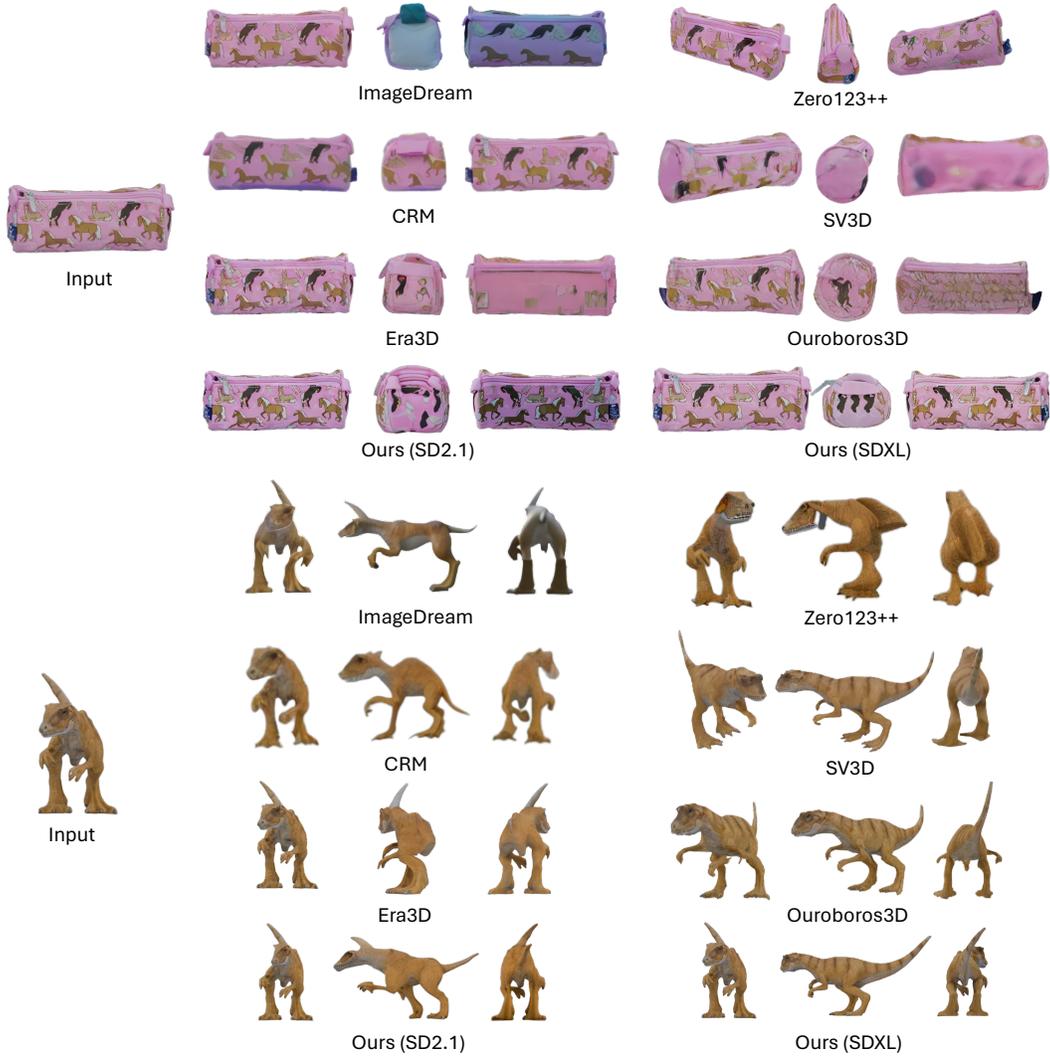


Figure 21: More qualitative comparison on image-to-multiview generation.

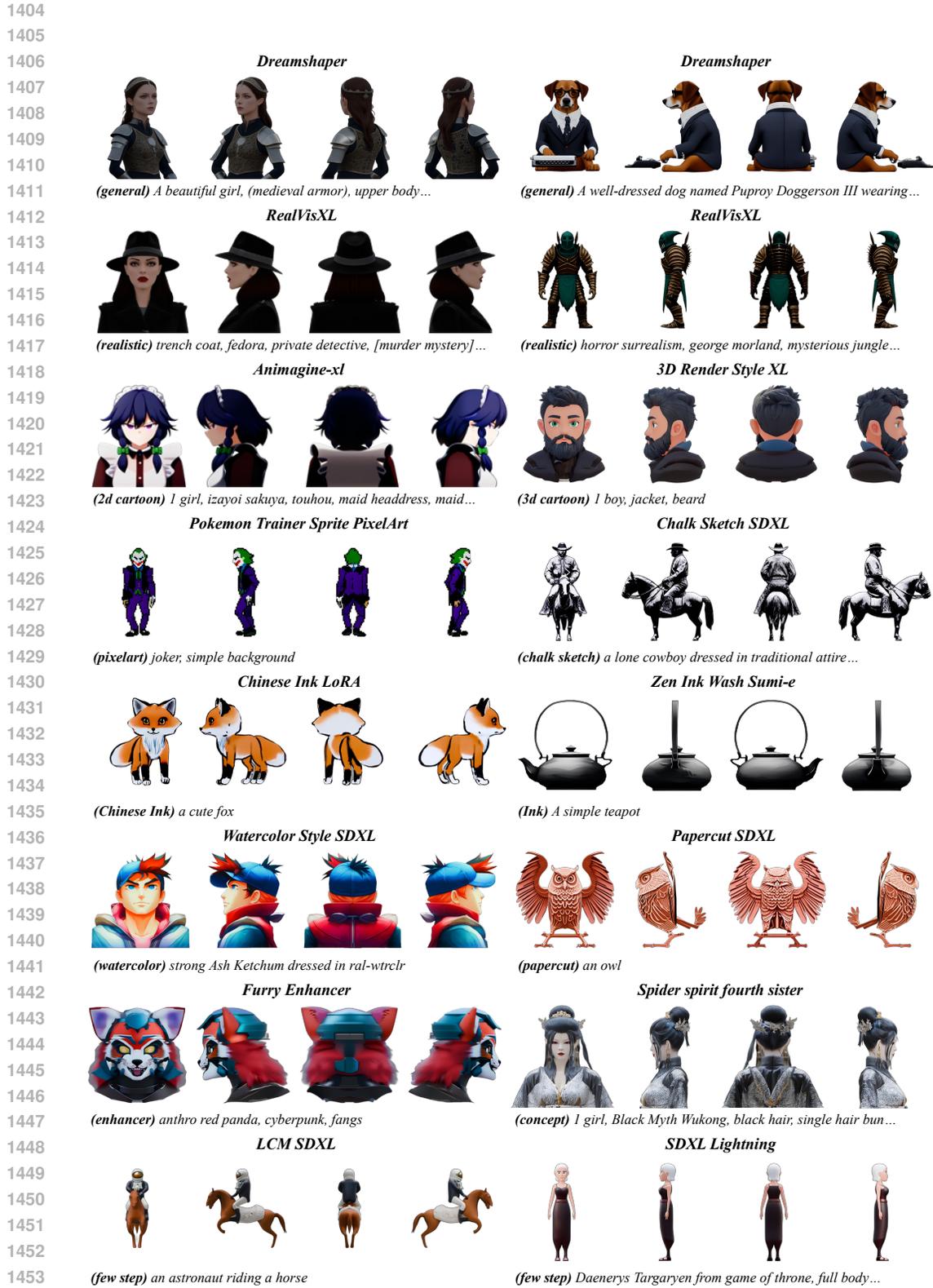


Figure 22: Additional results on camera-guided text-to-multiview generation with community models.

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

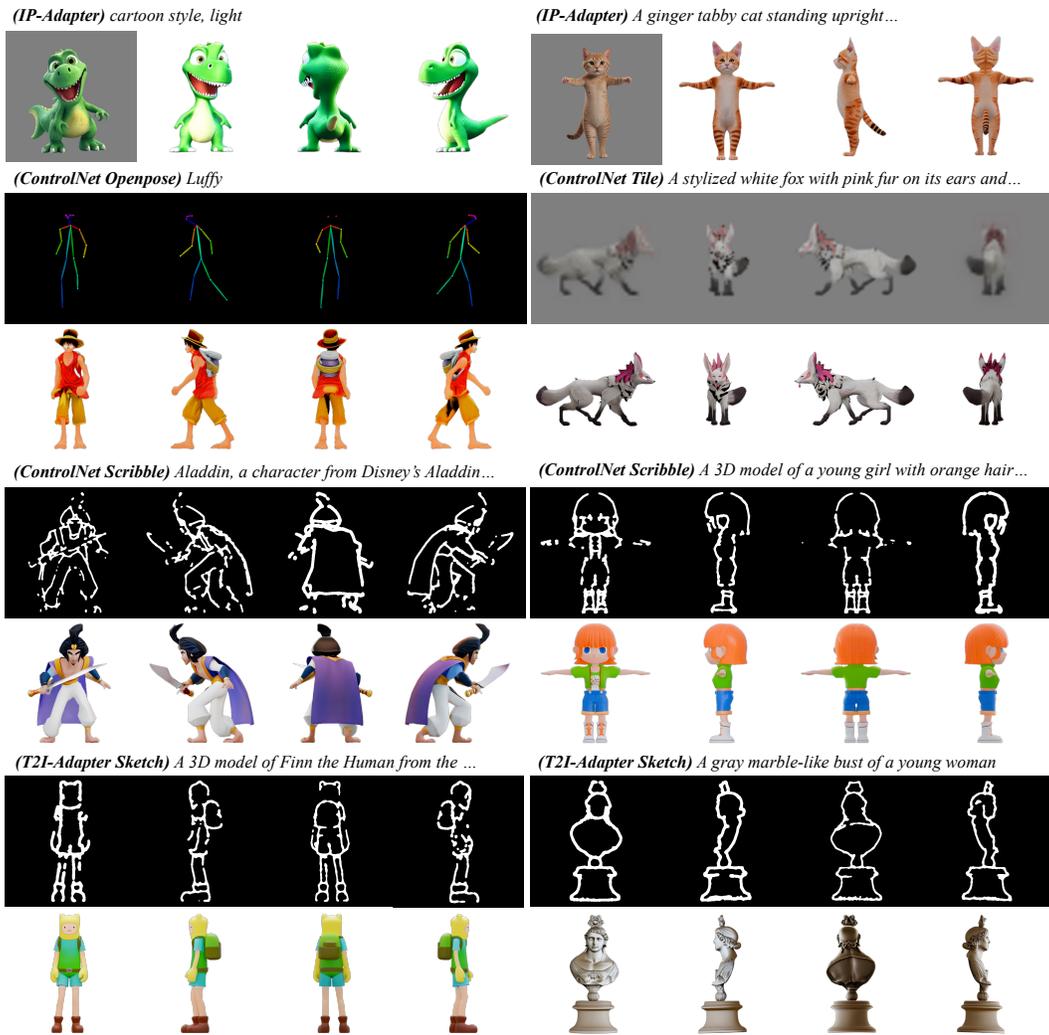


Figure 23: Additional results on camera-guided text-to-multiview generation with extensions.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

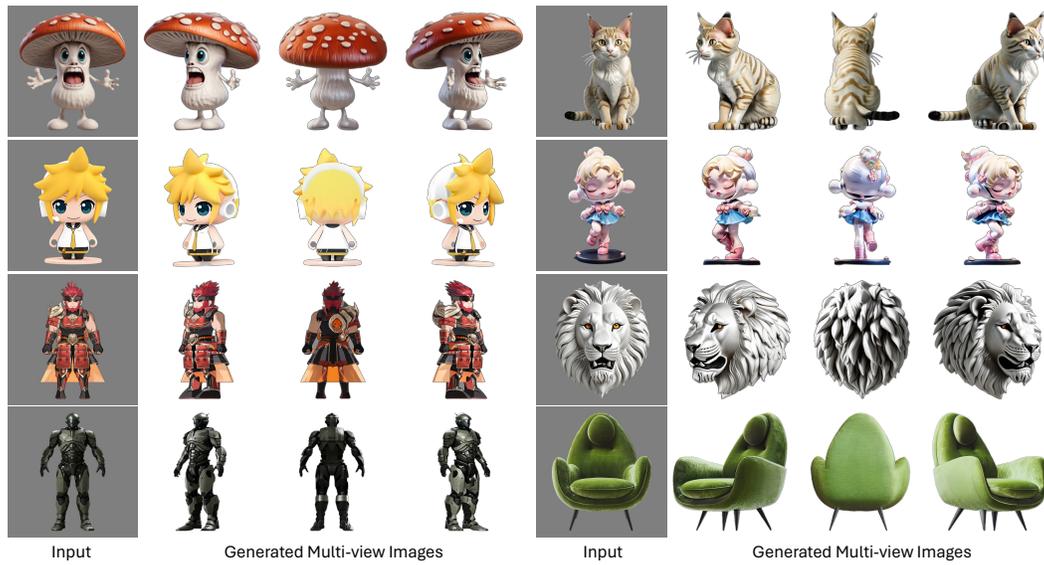


Figure 24: Additional results on camera-guided image-to-multiview generation.



Figure 25: Visual results on text-to-3D generation.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619



Figure 26: Visual results on image-to-3D generation.

*A blue low poly formula one car with number 33 on the body and a white circle on the hood...*



*A cartoon-styled rocket ship ride with a predominantly orange body, a green base, white details.*



*A US army motorcycle with a medical cross on the sidecar, a headlight, a brown seat, a large wheel...*



*Mater, a rusty and beat-up tow truck from the 2006 Disney/Pixar animated film "Cars", with a rusty...*



*A young girl with black hair, wearing an orange dress and yellow shirt, from the waist up.*



*A stylized squirrel holding an acorn, chubby body, short legs, a large fluffy tail, and big round eyes.*



*The 3D model is of the Super Sonic, a yellow anthropomorphic hedgehog with spiky hair...*



*A robot with blue, red and gray colors, and has a flame-like pattern on the body...*



*A purple anthropomorphic chameleon with a yellow belly and purple eyes, wearing black and purple...*



*Coco Bandicoot, from the Crash Bandicoot series, wearing her signature orange shirt and blue overalls...*



Figure 27: Additional results on geometry-guided text-to-texture generation.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

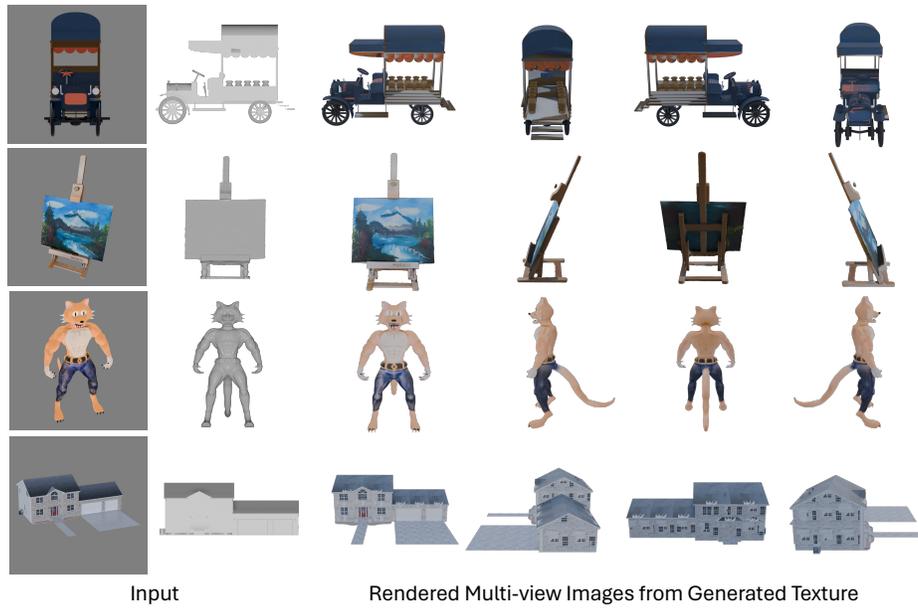


Figure 28: Additional results on geometry-guided image-to-texture generation.