# Scene123: One Prompt to 3D Scene Generation via Video-Assisted and Consistency-Enhanced MAE

**Yiying Yang[1*], Fukun Yin[2*], Jiayuan Fan[1†], Wanzhang Li[2], Xin Chen[3], Gang Yu[3]**

[1] Academy for Engineering and Technology, Fudan University
[2] School of Information Science and Technology, Fudan University
[3] Tencent PCG

yiyingyang23@m.fudan.edu.cn, fkyin21@m.fudan.edu.cn, jyfan@fudan.edu.cn,
liwz22@m.fudan.edu.cn, chenxin2@shanghaitech.edu.cn, iskicy@gmail.com

## Abstract

As Artificial Intelligence Generated Content (AIGC) advances, a variety of methods have been developed to generate text, images, videos, and 3D shapes from single or multimodal inputs, contributing efforts to emulate human-like cognitive content creation. However, generating realistic large-scale scenes from a single input presents a challenge due to the complexities involved in ensuring consistency across extrapolated views generated by models. Benefiting from recent video generation models and implicit neural representations, we propose Scene123, a 3D scene generation model, which combines a video generation framework to ensure realism and diversity with implicit neural fields integrated with Masked Autoencoders (MAE) to effectively ensure the consistency of unseen areas across views. Specifically, the input image (or a text-generated image) is first warped to simulate adjacent views, with the invisible regions filled using the consistency-enhanced MAE model. Nonetheless, the synthesized images often exhibit inconsistencies in viewpoint alignment, thus we utilize the produced views to optimize a neural radiance field, enhancing geometric consistency. Moreover, to further enhance the details and texture fidelity of generated views, we employ a GAN-based Loss against images derived from the input image through the video generation model. Extensive experiments demonstrate that our method can generate realistic and consistent scenes from a single prompt. Both qualitative and quantitative results indicate that our approach surpasses existing state-of-the-art methods. We show encourage video examples at https://yiyingyang12.github.io/Scene123.github.io/.

## 1 Introduction

3D scene generation aims to create realistic or stylistically specific scenes from limited prompts, such as a few images or a text description. This is a fundamental issue in computer vision and graphics and a critical challenge in generative artificial intelligence. Recent advancements have demonstrated substantial progress through the use of vision-language models (Radford et al. 2021), generative models such as Generative Adversarial Networks (GANs) (Zhang et al. 2019), Variational autoencoders (VAEs) (Sargent et al. 2023a; Yang et al. 2023b), or diffusion models (Yang et al. 2023a), and scene
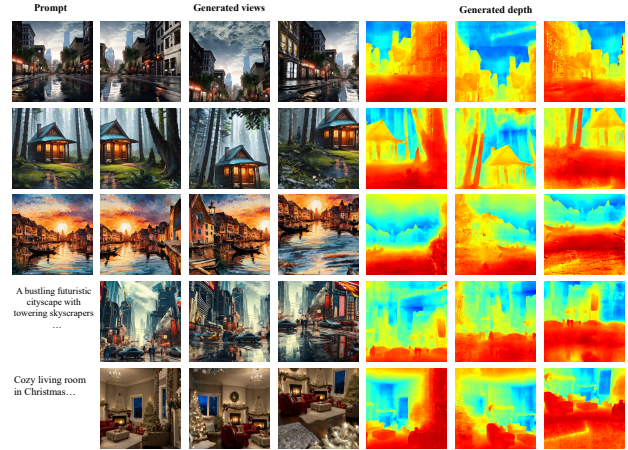


Figure 1: Some examples generated by our Scene123. For a single input image or text, our method can generate 3D scenes with consistent views, fine geometry, and realistic textures, applicable to real, virtual, or object-centered scenes.

representations like Neural Radiance Fields (NeRF) (Mildenhall et al. 2021; Yin et al. 2022; Ding et al. 2023; Lu et al. 2023; Yang et al. 2024) or 3D Gaussian Splatting (3DGS) (Kerbl et al. 2023). Typically, these methods begin by generating images with a pretrained generative model (Rombach et al. 2022), or by directly using images for image-to-scene generation, then estimate additional 3D surface geometric details such as depth and normal (Piccinelli, Sakaridis, and Yu 2023; Li et al. 2023; Yin et al. 2023; Yin and Zhou 2020), and subsequently render the surface textures of the scenes using generative strategies like inpainting (Wang et al. 2023b). Moreover, some approaches update the geometric surfaces as new views are generated to maintain the coherence of the scene (Zhang et al. 2024b). However, these methods often rely on pre-trained models (Radford et al. 2021; Li et al. 2022), resulting in inconsistencies and artifacts in the generated scenes. Additionally, these methods face challenges in producing high-quality, coherent 3D representations across diverse and complex environments.

In this paper, we endeavor to address this challenge about generating 3D scenes from a single image or textual description, ensuring viewpoint consistency and realistic surface tex-

---

[*]Yiying Yang and Fukun Yin contributed equally to this work.
[†]Corresponding author.

tures for both real and synthetically styled scenes, as shown in Fig. 1. However, achieving a balance between viewpoint consistency and flexibility presents a significant challenge. View consistency requires maintaining coherent and accurate details across multiple perspectives, which can limit the model's adaptability to diverse inputs and tasks. Conversely, flexibility requires the model to produce high-quality outputs under varying conditions, which may introduce inconsistencies in the generated scenes.

It is noteworthy that scene synthesis based on multi-view images is a long-studied topic from the early Multi-View Stereo (MVS) (Schönberger et al. 2016; Schönberger and Frahm 2016) to recent implicit neural representations (Yu et al. 2021), ensuring view consistency in generated scenes through multi-view matching mechanisms. However, when the input is reduced to a single image or a textual description, lacking reference multi-view images, the performance of these methods is greatly compromised. Fortunately, methodologies such as Masked Autoencoders (MAE) (He et al. 2022) provide avenues for extrapolating to areas unseen in new views, while the incorporation of additional semantic layers augments the coherence of the synthesized scenes. Additionally, video generation models (Luo et al. 2023) facilitate the enhancement of scenes with richer priors and more detailed texture information. Thus, how to utilize multi-view reconstruction with robust physical constraints alongside effective expansions in new viewpoints and video generation models is still an unreached area.

Scene123 investigates a methodology for 3D scene reconstruction employing stringent physical constraints alongside a robust multi-view MAE and video generation models to ensure view consistency. To achieve this, for each input or generated single image, we first create images of nearby perspectives through warping. Subsequently, a consistency-enhanced MAE is designed to inpaint the unseen areas, with a shared codebook maintained to distribute global information to every invisible area. We implement a progressive strategy derived from Text2NeRF (Zhang et al. 2024b) to incrementally update perspectives throughout this process. Furthermore, to enhance the scene's detail and realism, we employ the latest video generation technology, generating high-quality scene videos based on input images and performing adversarial enhancements with the rendered images. The synergistic operation of these two modules enables Scene123 to generate consistent, finely detailed three-dimensional scenes with photo-realistic textures from a single prompt.

We conduct extensive experiments on text-to-scene and image-to-scene generation, encompassing both real and virtual scenes. The results offer robust empirical evidence that strongly supports the effectiveness of our framework. Our contributions can be summarized as: 1) We propose a novel scene generation framework based on one prompt, which establishes the connection between MAE and video generation models for the first time to ensure view consistency and realism of the generated scenes. 2) The Consistency-Enhanced MAE is designed to fill unseen areas in new views by injecting global semantic information and combining it with neural implicit fields, ensuring consistent surface representation across various views. 3) We introduce the video-assisted

3D-aware generative refinement module, which enhances scene reconstruction by integrating the diversity and realism of video generation models through a GAN-based function to significantly improve detail and texture fidelity. 4) Extensive experiments validate the efficacy of Scene123, demonstrating greater accuracy in surface reconstruction, higher realism in reconstructed views, and better texture fidelity compared to the SOTA methods. The data and code will be available.

## 2 Related Work

**Text to 3D Scene Generation.** Many recent advancements in text-driven 3D scene generation have focused on modeling 3D scenes using text inputs (Hwang, Kim, and Kim 2023; Zhang et al. 2024b). Due to the scarcity of paired text-3D scene data, most studies utilize Contrastive Language-Image Pre-training (CLIP) (Radford et al. 2021) or pre-trained text-to-image models to interpret the text input. Text2Scene (Hwang, Kim, and Kim 2023) employs CLIP to model and stylize 3D scenes from text (or image) inputs by decomposing the scene into manageable sub-parts. Set-the-Scene (Cohen-Bar et al. 2023) and Text2NeRF (Zhang et al. 2024b) generate NeRFs from text using text-to-image diffusion models to represent 3D scenes. SceneScape (Fridman et al. 2024) and Text2Room (Höllein et al. 2023) leverage a pre-trained monocular depth prediction model for enhanced geometric consistency and directly generate the 3D textured mesh representation of the scene. However, these methods depend on inpainted images for scene completion, which, while producing realistic visuals, suffer from limited 3D consistency. More recent studies (Bai et al. 2023) have successfully generated multi-object compositional 3D scenes. Another class of methods utilizes auxiliary inputs, such as layouts (Po and Wetzstein 2023), to enhance scene generation. Unlike text captions, which can be rather vague, some approaches generate 3D scenes from image inputs, where the 3D scene corresponds closely to the depicted image. Early scene generation methods often require specific scene data for training to obtain category-specific scene generators, such as GAUDI (Bautista et al. 2022), or focus on single scene reconstruction based on the input image, such as PixelSynth (Rockwell, Fouhey, and Johnson 2021) and Worldsheet (Hu et al. 2021). However, these methods are often limited by the quality of the generation or the extensibility of the scene.

**Image to 3D Scene Generation.** Recently, numerous studies have concentrated on generating 3D scenes from image inputs. PERF (Wang et al. 2023a) generates 3D scenes from a single panoramic image, using diffusion models to supplement shadow areas. ZeroNVS (Sargent et al. 2023b) extends this capability by reconstructing both objects and environments in 3D from a single image. Despite its environmental reconstruction lacking some detail, the algorithm demonstrates an understanding of environmental contexts. Lucid-Dreamer (Chung et al. 2023) and WonderJourney (Yu et al. 2023) utilize general-purpose depth estimation models to project hallucinated 2D scene extensions into 3D representations. However, these methods still face challenges in achieving realism, often producing artifacts due to the reliance on pre-trained models.

**Video Diffusion Models and 3D-aware GANs.** Diffusion

models (Song et al. 2020) have recently emerged as powerful generative models capable of producing a wide array of images (Blattmann et al. 2023b) and videos (Blattmann et al. 2023a) by iteratively denoising a noise sample. Among these models, the publicly available Stable Diffusion (SD) (Rombach et al. 2022) and Stable Video Diffusion (SVD) (Blattmann et al. 2023a) exhibit strong generalization capabilities due to training on extremely large datasets such as LAION (Schuhmann et al. 2022) and LVD (Blattmann et al. 2023a). Consequently, they are frequently used as foundational models for various generation tasks, including novel view synthesis. To enhance generalization and multi-view consistency, some contemporary works leverage the temporal priors in video diffusion models for object-centric 3D generation. For instance, IM-3D (Melas-Kyriazi et al. 2024) and SV3D (Voleti et al. 2024) explore the capabilities of video diffusion models in object-centric multi-view generation. V3D (Chen et al. 2024b) extends this approach to scene-level novel view synthesis. However, these methods often produce unsatisfactory results for complex objects or scenes, leading to inconsistencies among multiple views or unrealistic geometries.

Early works, like 3D-GAN (Wu et al. 2016), Pointflow (Yang et al. 2019), and ShapeRF (Cai et al. 2020) focus more on the category-specific texture-less geometric shape generation based on the representations of voxels or point clouds. However, limited by the generation capabilities of GANs, these methods can only generate rough 3D assets of specific categories. Subsequently, HoloGAN (Nguyen-Phuoc et al. 2019), GET3D (Gao et al. 2022), and EG3D (Chan et al. 2022) employ GAN-based 3D generators conditioned on latent vectors to produce category-specific textured 3D assets. Recently, as seen in GigaGAN (Kang et al. 2023), the Generative Adversarial Networks (GANs) methods are better suited for high-frequency details than diffusion models. Furthermore, IT3D (Chen et al. 2024a) proposes a novel Diffusion-GAN dual training strategy to overcome the view inconsistency challenges. However, the training process of GAN is prone to the issue of mode collapse, which limits the diversity of generation results.

## 3 Methodology

### 3.1 Overview

In this paper, we propose a novel 3D scene generation framework based on one prompt, a single image or textual description, ensuring viewpoint consistency and realistic surface textures for both real and synthetically styled scenes, as shown in Fig. 2. We first design the Consistency-Enchanced MAE module to fill unseen areas in novel views by injecting global semantic information. With support views in the initialized database $\mathbf{S}$ generated by the Consistency-Enhanced MAE module, we employ a NeRF network to represent the 3D scene as the physical constraints of view consistency. Furthermore, to enhance the scene's detail and realism, we employ the latest video generation technology, generating high-quality scene videos based on input images and performing adversarial enhancements with the rendered images. Finally, the optimization and implementation details of the model are detailed.

### 3.2 Consistency-Enhanced MAE Scene Completion

**Scene Initialization.** Given the reference image $\mathbf{I_0}$, notably, for only text prompt input, we utilize the stable diffusion model (Rombach et al. 2022) to generate initial image $\mathbf{I_0}$, we then feed this image into the off-the-shelf depth estimation model (Miangoleh et al. 2021), and take the output as a geometric prior for the target scene, denoted as $\mathbf{D_0}$. Inspired by (Zhang et al. 2024b), we construct an original database $\mathcal{S}_0 = \{(\mathbf{I}_i, \mathbf{D}_i)\}_{i=1}^N$ via the depth image-based rendering (DIBR) method (Fehn 2004), where $N$ denotes the number of initial viewpoints. Specifically, for each pixel $x$ in $\mathbf{I_0}$ and its depth value $y$ in $\mathbf{D_0}$, we compute its corresponding pixel $x_{0 \to m}$ and depth $y_{0 \to m}$ on a surrounding view $m$, $[x_{0 \to m}, y_{0 \to m}]^T = \mathbf{K}\mathbf{P}_m\mathbf{P}_0^{-1}\mathbf{K}^{-1}[x, y]^T$, where $\mathbf{K}$ and $\mathbf{P}_m$ indicate the intrinsic matrix and the camera pose in view $m$. This database provides additional views and depth information, which could prevent the model from overfitting to the initial view.

**MAE Scene Completion.** However, the original database $\mathcal{S}_0$ will inevitably have missing content since the information in the initial scene is derived from the single image $\mathbf{I_0}$ to construct the initialized database $\mathcal{S}$. Directly applying the original database to the 3D scene representation would inevitably suffer from limited 3D consistency of the generated scenes (Zhang et al. 2024b; Höllein et al. 2023). Inspired by (He et al. 2022; Hu et al. 2023; Zhang et al. 2024a), we design a Consistency-Enhanced MAE module to effectively ensure the consistency of unseen areas across views. Specifically, we first design a discrete codebook to distribute global information to every invisible area in the original database. We represent the codebook as $\mathcal{E} = \{e_1, e_2, ..., e_N\} \in \mathbb{R}^{N \times n_q}$, where $N$ stands for the total count of prototype vectors, $n_q$ denotes the dimensionality of individual vectors, and $e_i$ symbolizes each specific embedding vector. Specifically, given the input image $x \in \mathbb{R}^{H \times W \times 3}$, the VQ-VAE (Van Den Oord, Vinyals et al. 2017) employs an encoder $E$ to extract a continuous feature representation: $\hat{z} = E(x) \in \mathbb{R}^{h \times w \times c}$, where $h$ and $w$ are the height and width of the feature map. This continuous feature map $\hat{z}$ is then subjected to a quantization process $Q$, aligning it with its nearest codebook entry $e_k$ to obtain its discrete representation $z_q$ as follows:

$$z_q = Q(\hat{z}) := \operatorname*{argmin}_{e_k \in \mathcal{E}} \|\hat{z}_{ij} - e_k\|_2, \tag{1}$$

where $\hat{z}_{ij} \in \mathbb{R}^c$. We pre-trained VQVAE on the ImageNet (Russakovsky et al. 2015) dataset to equip the codebook with a more representative and generalized feature set. Subsequently, we fine-tuned it on specific scenes to better capture and represent the unique characteristics of each scene. Moreover, we then utilize the MAE encoder to encode the input images from the original database $\mathcal{S}_0 = \{(\mathbf{I}_i, \mathbf{D}_i)\}_{i=1}^N$, forming the image-conditional $\mathbf{s_c} = \{s_1, s_2, ..., s_M\}$, each embedding vector $s_i$ queries the valuable prior information from the given codebook via the cross-attention mechanism,

$$Q \leftarrow f_q(\mathbf{s_c}), \qquad K \leftarrow f_k(\mathcal{E}), \qquad V \leftarrow f_v(\mathcal{E})$$

$$\mathbf{s_u} \leftarrow \text{Cross-Attention}(Q, K, V) = \text{Softmax}(\frac{QK^T}{\sqrt{d_k}}), \tag{2}$$
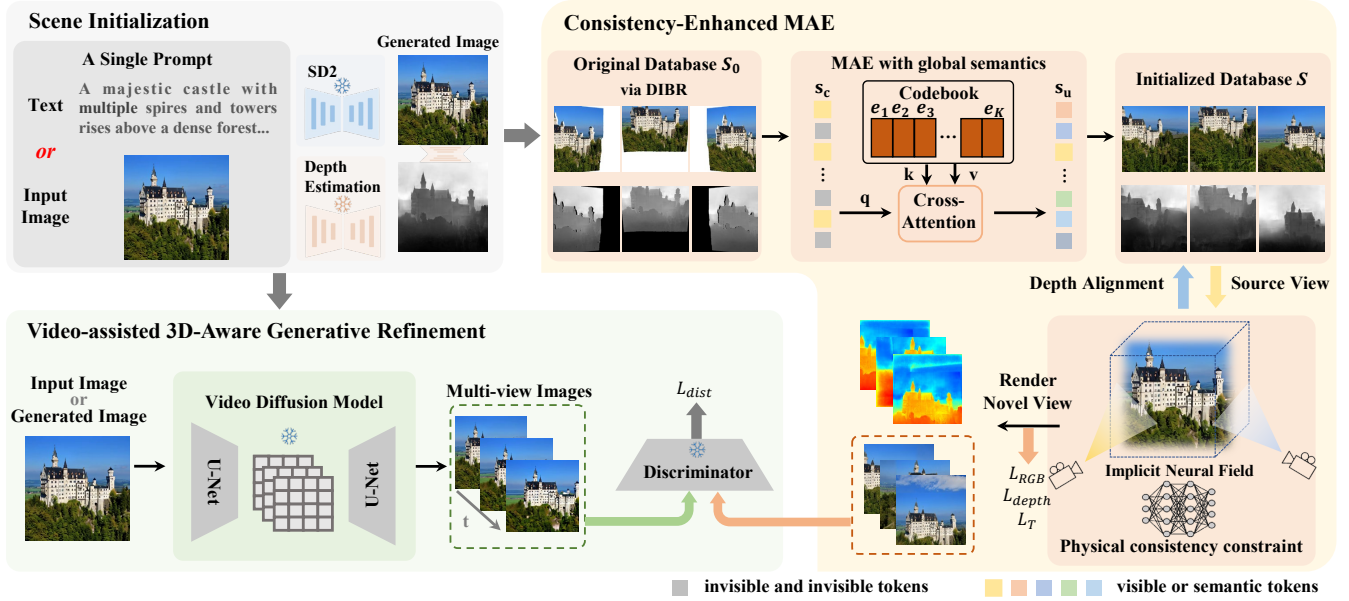
Figure 2: Scene123's pipeline includes two key modules: the consistency-enhanced MAE and the 3D-aware generative refinement module. The former generates adjacent views from an input image via warping, using the MAE model to inpaint unseen areas with global semantics and optimizing an implicit neural field for viewpoint consistency. The latter generates realistic videos from the input image with a pre-trained video generation model, enhancing realism through adversarial loss with rendered images.

where $f_q$, $f_k$, and $f_v$ are the query, key, and value linear projections, respectively. Consequently, the global semantic information contained in the codebook is maintained to distribute global information to every invisible area. Then, we utilize the MAE decoder to decode the feature $s_u$, deriving the initialized database $\mathcal{S}$, as shown in Fig. 2. Using global semantics as Key and Value in cross-attention allows the model to integrate comprehensive scene information, maintaining coherence. This approach enhances multi-view consistency by providing a unified understanding of the scene, reducing artifacts, and ensuring seamless integration of novel views.

### 3.3 3D Scene Representation

With these support views in the initialized database $\mathcal{S}$, along with the initial view $\mathbf{I_0}$, we aim to generate 3D scenes through robust 3D scene representations to provide physical-level surface consistency constraints. In this work, we employ a NeRF network $f_\theta$ to represent the 3D scene. Specifically, in NeRF representation, volume rendering (Mildenhall et al. 2021) is used to accumulate the color in the radiance fields,

$$\mathbf{C}(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})\mathrm{d}t, \qquad (3)$$

where $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ represents the 3D coordinates of sampled points on the camera ray emitted from the camera center $\mathbf{o}$ with the direction $\mathbf{d}$. $t_n$ and $t_f$ indicate the near and far sampling bounds. $(\mathbf{c}, \sigma) = f_\theta(\mathbf{r}(t))$ are the predicted color and density of the sampled point along the ray.

$$T(t) = \exp\left(-\int_{t_n}^{t} \sigma(\mathbf{r}(s))\mathrm{d}s\right), \qquad (4)$$

where $T(t)$ is the accumulated transmittance. Different from NeRF that takes both the 3D coordinate $\mathbf{r}(t)$ and view direction $\mathbf{d}$ to predict the radiance $\mathbf{c}(\mathbf{r}(t), \mathbf{d})$, we omit $\mathbf{d}$ to avoid the effect of view-dependent specularity.

However, due to the lack of geometric constraint during the depth estimation, the predicted depth values could be misaligned in the overlapping regions (Luo et al. 2020). Despite the design of the consistency-enhanced MAE module greatly improving the inconsistency between different views, the estimated depth rendered from NeRF still may be inconsistent between views. Inspired by Text2NeRF (Zhang et al. 2024b), we globally align these two depth maps by compensating for mean scale and value differences. Specifically, we first we first perform global alignment by calculating the average $s$ and depth offset $\delta$ to approximate the mean scale and value differences, and then we finetune a pre-trained depth alignment network to produce a locally aligned depth map. More details will be shown in the Appendix.

### 3.4 Video-assisted 3D-Aware Generative Refinement

**Support Set Generation.** Given the Reference image $\mathbf{I_0}$, we employ a image-to-video pipeline to generate a support set of enhanced quality, termed $D$, which is conditioned on the input reference image. For the image-to-video pipeline, we opt for Stable Video Diffusion(SVD) (Blattmann et al. 2023a), which is trained to generate smooth and consistent videos on large-scale datasets of real and high-quality videos. The exposure to superior data quantity and quality makes it more generalizable and multi-view consistent, and the flexibility of the SVD architecture makes it amenable to be fine-tuned for camera controllability. In the context of the SVD
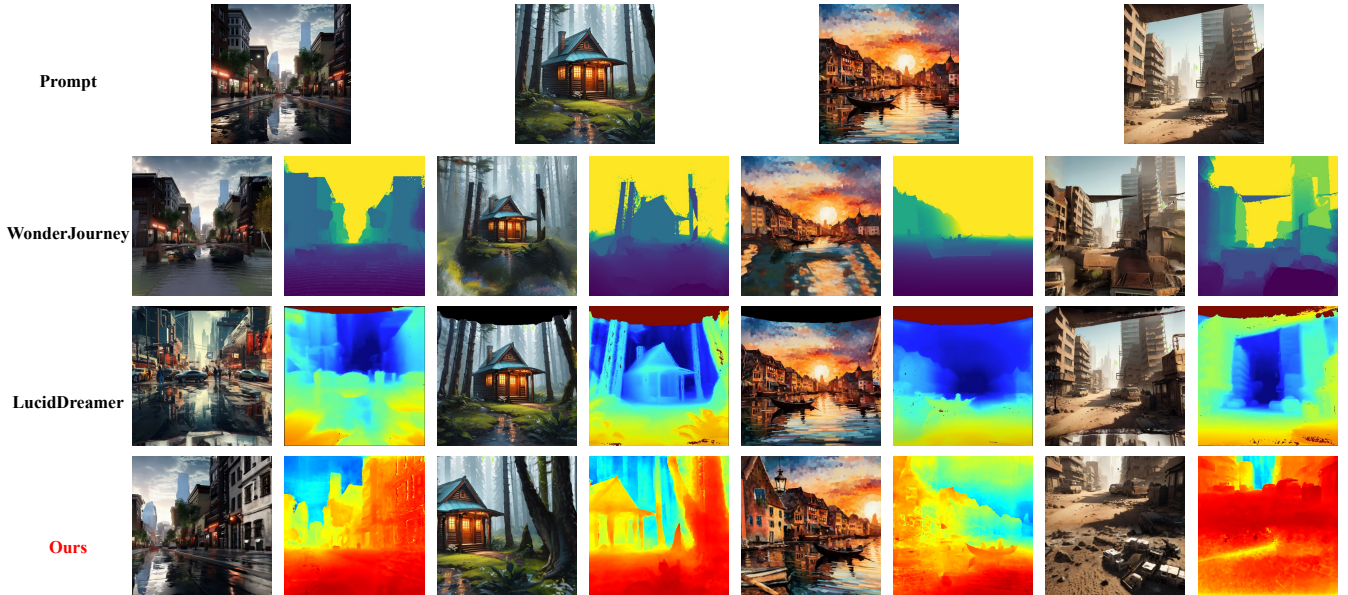
Figure 3: Qualitative results (zoom-in to view better) of methods capable of processing a single image prompt. We both visualize the texture and depth from novel views within the scene.

image-to-video (Blattmann et al. 2023a) pipeline, a noise-augmented (Ho et al. 2022) version of the conditioning frame channel-wise is concatenated to the input of the UNet (Ronneberger, Fischer, and Brox 2015). In addition, the temporal attention layers in SVD naturally assist in the consistent multi-view generation without needing any explicit 3D structures like in (Liu et al. 2023).

**3D-Aware Generative Refinement.** The capabilities of Generative Adversarial Networks (GANs) shine in scenarios involving datasets characterized by high variance (Chan et al. 2022). GANs have the ability to learn both geometry and texture-related knowledge from datasets, subsequently guiding the model to converge towards the same high-quality distribution exhibited by the generated support set. In our approach, we designate the video diffusion model as a generator. As shown in Fig. 4, given a reference image, SVD can generate consistent multi-view images at one time, constructing the generated support set. We then incorporate a discriminator initialized with random values. In this setup, the Support Set $D$ is treated as real data, while the renderings of the 3D neural radiance field model represent fake data. The role of the discriminator involves learning the distribution discrepancy between the renderings and $D$, subsequently contributing to the discrimination loss, which in turn updates the 3D neural radiance field model. This 3D-aware generative refinement module utilizes the discrimination loss $L_{\mathbf{dist}}$, which can help guide the updating direction and enhance the model's ability to produce intricate geometry and texture details.

The camera motion in the SVD model is sometimes limited(Blattmann et al. 2023a), making it unsuitable for directly training the NeRF model. Instead, we use our 3D-Aware Generative Refinement approach, which leverages a discriminator to optimize the process. Despite the small range of camera motion in the images generated by SVD, the adversarial train-



Figure 4: Data samples generated via image-to-video model.

ing process provides necessary regularization, demonstrating that even with limited camera motion, the discrimination loss remains effective and contributes to generating high-quality 3D scenes.

## 3.5 Optimization and Implementation Details

In addition to discrimination loss, we also utilize a $L_2$ loss, depth loss, and a transmittance loss to optimize the radiance field of the 3D scene, following previous NeRF-based works (Chen et al. 2022; Song et al. 2023; Zhang et al. 2024b). The RGB loss $L_{\mathbf{RGB}}$ is defined as a $L_2$ loss between the render pixel $\boldsymbol{C}^R$ and the color $\boldsymbol{C}$ generated by the MAE model. Different from previous works that employ regularized depth losses to handle uncertainty or scale-variant problems (Roessle et al. 2022; Sargent et al. 2023a), we adopt a stricter depth loss $L_{\mathbf{depth}}$ to minimize the $L_2$ distance between the rendered depth and the estimated depth. Moreover, we compute a depth-aware transmittance loss $L_{\mathbf{T}}$ (Jain et al. 2022; Zhang et al. 2024b) to encourage the NeRF network to produce empty density before the camera ray reaches the expected depth $\hat{\mathbf{z}}$, $L_{\mathbf{T}} = \|\mathbf{T}(t) \cdot \mathbf{m}(t)\|_2$, where $\mathbf{m}(t)$ is a mask indicator that satisfies $\mathbf{m}(t) = 1$ when $t < \hat{\mathbf{z}}$, otherwise $\mathbf{m}(t) = 0$. $\hat{\mathbf{z}}$ is the pixel-wise depth value in the estimated depth map, and $\mathbf{T}(t)$ is the accumulated transmittance. Therefore, the total loss function is then defined as,
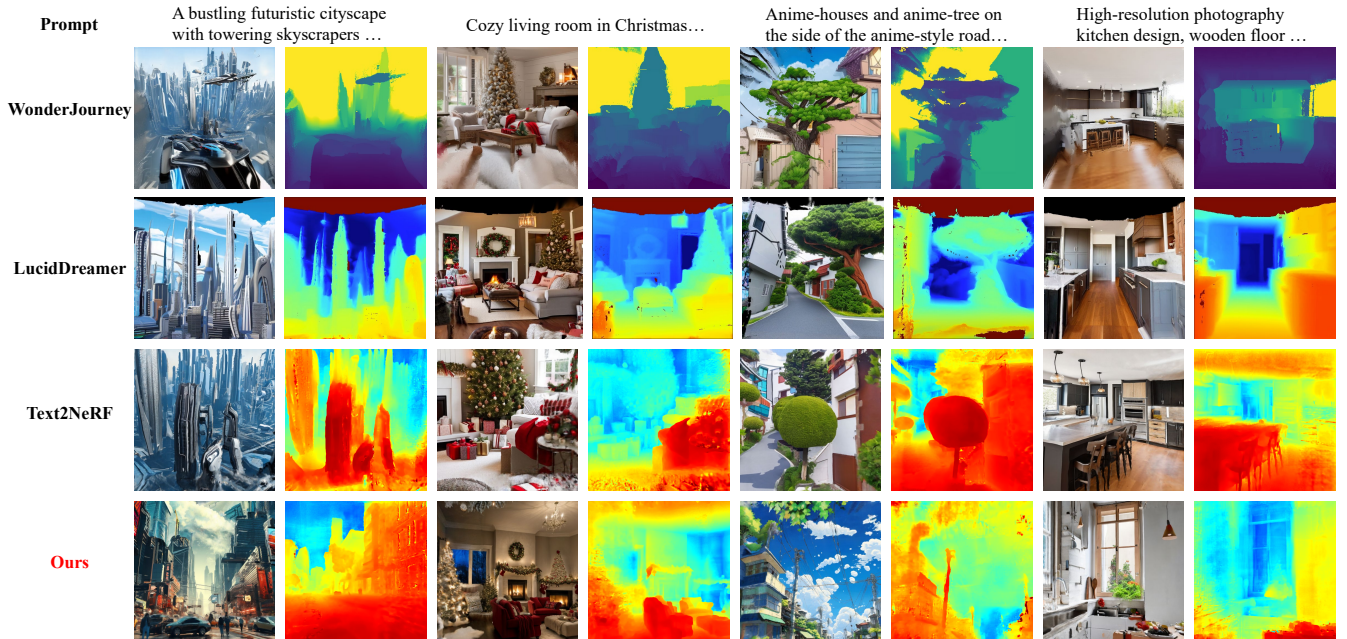
Figure 5: Qualitative results (zoom-in to view better) of methods that generate scenes from textual input. We both visualize the texture and depth from novel views within the scene.

$$L_{\mathbf{total}} = L_{\mathbf{RGB}} + \lambda_d L_{\mathbf{depth}} + \lambda_t L_{\mathbf{T}} + \lambda_{dist} L_{\mathbf{dist}}, \quad (5)$$

where $\lambda_d$, $\lambda_T$, $\lambda_{dist}$ are constant hyperparameters balancing depth, transmittance, and discrimination losses.

**Implementation Details.** We implement our Scene123 with the Pytorch framework (Paszke et al. 2019) and adopt TensoRF (Chen et al. 2022) as the radiance field. To ensure TensoRF can accommodate scene generation over a large view range, we position the camera near the center of the NeRF bounding box and configure it with outward-facing viewpoints. The dimension of the masked codebook is 2048×16. For only text prompt input, we utilize the stable diffusion model in version 2.0 (Rombach et al. 2022) to generated initial image $\mathbf{I_0}$ related to the input prompt. Moreover, for depth estimation, we use the boosting monocular depth estimation method (Miangoleh et al. 2021) with pre-trained LeReS model (Yin et al. 2021) to estimate the depth for each view. For the image-to-video pipeline, we opt for the Stable Video Diffusion in version SVD-XT, which is the same architecture as SVD (Blattmann et al. 2023a) but finetuned for 25 frame generation. During training, we use the same setting as (Chen et al. 2022) for the optimizer and learning rate and set the hyperparameters in our objective function as $\lambda_d = 0.005$, $\lambda_t = 0.001$, $\lambda_{dist} = 0.001$.

# 4 Experiments

## 4.1 Experimental Setup

**Dataset and baselines.** Since the perpetual 3D scene generation is a new task without an existing dataset, we use real or generated high-quality images as input (Yu et al. 2023) for evaluation in our experiments. We consider two state-of-the-art 3D scene generation methods as our baselines,

WonderJourney (Yu et al. 2023) and LucidDreamer (Chung et al. 2023) to compare the performance of the image or text prompt as a condition input for 3D scene generation. WonderJourney and LucidDreamer rely on an off-the-shelf general-purpose depth estimation model to project the hallucinated 2D scene extensions into a 3D representation. Specifically, WonderJourney designs a fully modularized model to generate sequences of 3D scenes. LucidDreamer utilizes Stable Diffusion (Rombach et al. 2022) and 3D Gaussian splatting (Kerbl et al. 2023) to create diverse high-quality 3D scenes. For text prompts, besides WonderJourney and LucidDreamer, we include Text2NeRF (Zhang et al. 2024b) as the baseline, which performs well for the text-to-3D generation. Text2NeRF generates NeRFs from text with the aid of text-to-image diffusion models to represent 3D scenes.

**Evaluation Metrics.** Following Text2NeRF, we evaluate the quality of our generated images using CLIP Score (CLIP-Similarity), Blind/Referenceless Image Spatial Quality Evaluator (BRISQUE) (Mittal, Moorthy, and Bovik 2012) and Natural Image Quality Evaluator (NQIE) (Mittal, Soundararajan, and Bovik 2012).

## 4.2 Performance Comparisons

As shown in Tab. 1, we evaluate the quality of generated 3D scenes across baselines quantitatively and report the average evaluation scores of CLIP-Similarity BRISQUE and NIQE for the generated images produced by different methods. Clearly, our method surpasses the baselines by generating higher-quality 3D scenes, as indicated by lower BRISQUE and NIQE values. Moreover, our method ensures the semantic relevance between the generated scene and the input text, resulting in a higher CLIP score. The qualitative results are drawn in Fig. 3 and Fig. 5. Our method can ensure more viewpoint consistency and realistic surface textures for both

| Method | Input | Visual Quality | | | Optimization Time (hours) |
|---|---|---|---|---|---|
| | | CLIP-Similarity↑ | BRISQUE↓ | NIQE↓ | |
| WonderJourney (Yu et al. 2023) | Image&Text | 27.480 | 67.012 | 12.022 | **0.208** |
| LucidDreamer (Chung et al. 2023) | Image&Text | 26.663 | 46.266 | 6.652 | 0.220 |
| Text2NeRF (Zhang et al. 2024b) | Text | 28.695 | 24.498 | 4.618 | 1.525 |
| **Scene123 (Ours)** | Image&Text | **30.544** | **20.324** | **2.522** | 1.433 |

Table 1: Quantitative comparison results of our method with the baseline WonderJourney, LucidDreamer and Text2NeRF ↑ means the higher, the better, ↓ means the lower, the better.

real and synthetically styled scenes, giving a single image or textual description. Obviously, Fig. 3, given a reference image, WonderJourney and LucidDreamer seem to generate 3D scenes with artifacts, while our method can generate more semantic relevant to the given image, maintaining the multi-view consistency. In Fig. 5, with textual prompt input, WonderJourney can inevitably generate artifacts. While LucidDreamer and Text2NeRF are likely to produce distorted or inconsistent viewpoint images. Consequently, our method demonstrates superior qualitative performance compared to these baseline approaches.

## 4.3 Ablation Studies and Analysis

**Effectiveness of the Consistency-Enhanced MAE.** To verify the effectiveness of the Consistency-Enhanced MAE, we conduct ablation studies on different strategies: removing the MAE scene completion, denoted as w/o MAE; removing the masked VQ-VQE codebook, denoted w/o codebook; removing both MAE scene completion and the masked VQ-VAE codebook, denoted as w/o MAE&codebook. For w/o MAE, we utilize the stable diffusion (sd) inpainting to complete the scene, which is utilized in LucidDreamer. As shown in Tab. 2 and Fig. 6, the integration of the consistency-enhanced MAE and the codebook significantly contributes to the generation of coherent 3D scenes. Fig. 6 demonstrates visually that when the model removes either of the MAE and codebook modules, it causes inconsistencies between the different views. This superiority of our Consistency-Enhanced MAE in handling detailed complementation is evident in Fig. 6.

| Method | CLIP↑ | BRISQUE↓ | NIQE↓ |
|---|---|---|---|
| w/o MAE | 27.745 | 41.243 | 6.024 |
| w/o codebook | 27.983 | 38.335 | 5.834 |
| w/o MAE&codebook | 26.674 | 44.234 | 6.653 |
| **full model** | **30.544** | **20.324** | **2.522** |

Table 2: Ablation experiments regarding the MAE module and the codebook.



Figure 6: Results of the effectiveness of the MAE module and the codebook.

**Effectiveness of the 3D-aware generative refinement mod-**

**ule.** To verify the effectiveness of the video-assisted 3D-aware generative refinement module, we conduct ablation studies on different strategies, including removing the discrimination loss, denoted as w/o GAN loss; replacing the real support set generated by the image-to-video pipeline with a set containing the same number of images using reference image duplication, denoted as w/o video-assisted. As shown in Tab. 3 and Fig. 7, the incorporation of the GAN-based training strategy significantly enhances the model's ability to render detailed textures and complex geometries. The full model achieves the lowest FID values, while similar in other metrics, indicating the 3D-aware generative refinement module plays an important role in providing intricate geometry and texture details.

| Method | CLIP↑ | BRISQUE↓ | NIQE↓ |
|---|---|---|---|
| w/o GAN loss | 25.234 | 33.234 | 7.234 |
| w/o video-assisted | 28.341 | 24.342 | 5.342 |
| **full model** | **30.544** | **20.324** | **2.522** |

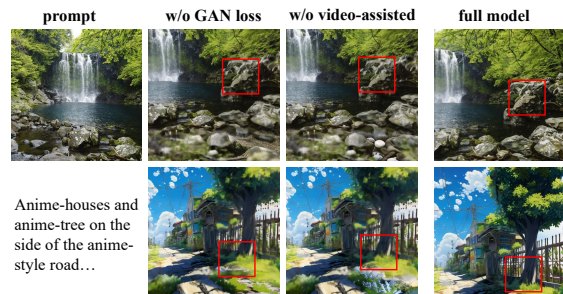Table 3: Ablation experiments regarding the 3D-aware generative refinement module.



Figure 7: Results of the effectiveness of the 3D-aware generative refinement module.

## 5 Conclusions and limitations

**Conclusions.** In this paper, we introduce Scene123, which surpasses existing 3D scene generation methods in scene consistency and realism, providing finer geometry and high-fidelity textures. Our method mainly relies on the consistency-enhanced MAE and the 3D-aware generative refinement module. The former exploits the inherent scene consistency constraints of implicit neural fields and integrates them with the MAE model with global semantics to inpaint adjacent views, ensuring viewpoint consistency. The latter uses the video generation model to produce realistic videos, enhancing the detail and realism of rendered views. With the help of

these two modules, our method can generate high-quality 3D scenes from a single input prompt, whether real, virtual, or object-centered settings.

**Limitations.** Our method is both innovative and effective. However, the optimization process remains an area for potential enhancement, as it is currently constrained by our use of implicit neural fields to enforce physical consistency. In our future research, we aim to explore faster, more easily optimized, and sustainable 3D scene representations. These will serve as mechanisms to articulate consistency surface constraints, thus accelerating the generation process.

# References

Bai, H.; Lyu, Y.; Jiang, L.; Li, S.; Lu, H.; Lin, X.; and Wang, L. 2023. Componerf: Text-guided multi-object compositional nerf with editable 3d scene layout. *arXiv preprint arXiv:2303.13843*.

Bautista, M. A.; Guo, P.; Abnar, S.; Talbott, W.; Toshev, A.; Chen, Z.; Dinh, L.; Zhai, S.; Goh, H.; Ulbricht, D.; et al. 2022. Gaudi: A neural architect for immersive 3d scene generation. *Advances in Neural Information Processing Systems*, 35: 25102–25116.

Blattmann, A.; Dockhorn, T.; Kulal, S.; Mendelevitch, D.; Kilian, M.; Lorenz, D.; Levi, Y.; English, Z.; Voleti, V.; Letts, A.; et al. 2023a. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*.

Blattmann, A.; Rombach, R.; Ling, H.; Dockhorn, T.; Kim, S. W.; Fidler, S.; and Kreis, K. 2023b. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 22563–22575.

Cai, R.; Yang, G.; Averbuch-Elor, H.; Hao, Z.; Belongie, S.; Snavely, N.; and Hariharan, B. 2020. Learning gradient fields for shape generation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, 364–381. Springer.

Chan, E. R.; Lin, C. Z.; Chan, M. A.; Nagano, K.; Pan, B.; De Mello, S.; Gallo, O.; Guibas, L. J.; Tremblay, J.; Khamis, S.; et al. 2022. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 16123–16133.

Chen, A.; Xu, Z.; Geiger, A.; Yu, J.; and Su, H. 2022. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, 333–350. Springer.

Chen, Y.; Zhang, C.; Yang, X.; Cai, Z.; Yu, G.; Yang, L.; and Lin, G. 2024a. It3d: Improved text-to-3d generation with explicit view synthesis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 1237–1244.

Chen, Z.; Wang, Y.; Wang, F.; Wang, Z.; and Liu, H. 2024b. V3d: Video diffusion models are effective 3d generators. *arXiv preprint arXiv:2403.06738*.

Chung, J.; Lee, S.; Nam, H.; Lee, J.; and Lee, K. M. 2023. Luciddreamer: Domain-free generation of 3d gaussian splatting scenes. *arXiv preprint arXiv:2311.13384*.

Cohen-Bar, D.; Richardson, E.; Metzer, G.; Giryes, R.; and Cohen-Or, D. 2023. Set-the-scene: Global-local training for generating controllable nerf scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2920–2929.

Ding, Y.; Yin, F.; Fan, J.; Li, H.; Chen, X.; Liu, W.; Lu, C.; YU, G.; and Chen, T. 2023. PDF: Point Diffusion Implicit Function for Large-scale Scene Neural Representation. arXiv:2311.01773.

Fehn, C. 2004. Depth-image-based rendering (DIBR), compression, and transmission for a new approach on 3D-TV. In *Stereoscopic displays and virtual reality systems XI*, volume 5291, 93–104. SPIE.

Fridman, R.; Abecasis, A.; Kasten, Y.; and Dekel, T. 2024. Scenescape: Text-driven consistent scene generation. *Advances in Neural Information Processing Systems*, 36.

Gao, J.; Shen, T.; Wang, Z.; Chen, W.; Yin, K.; Li, D.; Litany, O.; Gojcic, Z.; and Fidler, S. 2022. Get3d: A generative model of high quality 3d textured shapes learned from images. *Advances In Neural Information Processing Systems*, 35: 31841–31854.

He, K.; Chen, X.; Xie, S.; Li, Y.; Dollár, P.; and Girshick, R. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 16000–16009.

Ho, J.; Saharia, C.; Chan, W.; Fleet, D. J.; Norouzi, M.; and Salimans, T. 2022. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47): 1–33.

Höllein, L.; Cao, A.; Owens, A.; Johnson, J.; and Nießner, M. 2023. Text2room: Extracting textured 3d meshes from 2d text-to-image models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7909–7920.

Hu, Q.; Zhang, G.; Qin, Z.; Cai, Y.; Yu, G.; and Li, G. Y. 2023. Robust semantic communications with masked VQ-VAE enabled codebook. *IEEE Transactions on Wireless Communications*.

Hu, R.; Ravi, N.; Berg, A. C.; and Pathak, D. 2021. Worldsheet: Wrapping the world in a 3d sheet for view synthesis from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 12528–12537.

Hwang, I.; Kim, H.; and Kim, Y. M. 2023. Text2scene: Text-driven indoor scene stylization with part-aware details. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1890–1899.

Jain, A.; Mildenhall, B.; Barron, J. T.; Abbeel, P.; and Poole, B. 2022. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 867–876.

Kang, M.; Zhu, J.-Y.; Zhang, R.; Park, J.; Shechtman, E.; Paris, S.; and Park, T. 2023. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10124–10134.

Kerbl, B.; Kopanas, G.; Leimkühler, T.; and Drettakis, G. 2023. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4): 1–14.

Li, J.; Li, D.; Xiong, C.; and Hoi, S. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, 12888–12900. PMLR.

Li, S.; Zhou, J.; Ma, B.; Liu, Y.-S.; and Han, Z. 2023. Neaf: Learning neural angle fields for point normal estimation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, 1396–1404.

Liu, Y.; Lin, C.; Zeng, Z.; Long, X.; Liu, L.; Komura, T.; and Wang, W. 2023. Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*.

Lu, C.; Yin, F.; Chen, X.; Chen, T.; Yu, G.; and Fan, J. 2023. A Large-Scale Outdoor Multi-modal Dataset and Benchmark for Novel View Synthesis and Implicit Scene Reconstruction. *arXiv preprint arXiv:2301.06782*.

Luo, X.; Huang, J.-B.; Szeliski, R.; Matzen, K.; and Kopf, J. 2020. Consistent video depth estimation. *ACM Transactions on Graphics (ToG)*, 39(4): 71–1.

Luo, Z.; Chen, D.; Zhang, Y.; Huang, Y.; Wang, L.; Shen, Y.; Zhao, D.; Zhou, J.; and Tan, T. 2023. Videofusion: Decomposed diffusion models for high-quality video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10209–10218.

Melas-Kyriazi, L.; Laina, I.; Rupprecht, C.; Neverova, N.; Vedaldi, A.; Gafni, O.; and Kokkinos, F. 2024. IM-3D: Iterative Multiview Diffusion and Reconstruction for High-Quality 3D Generation. *arXiv preprint arXiv:2402.08682*.

Miangoleh, S. M. H.; Dille, S.; Mai, L.; Paris, S.; and Aksoy, Y. 2021. Boosting monocular depth estimation models to high-resolution via content-adaptive multi-resolution merging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9685–9694.

Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1): 99–106.

Mittal, A.; Moorthy, A. K.; and Bovik, A. C. 2012. No-reference image quality assessment in the spatial domain. *IEEE Transactions on image processing*, 21(12): 4695–4708.

Mittal, A.; Soundararajan, R.; and Bovik, A. C. 2012. Making a "completely blind" image quality analyzer. *IEEE Signal processing letters*, 20(3): 209–212.

Nguyen-Phuoc, T.; Li, C.; Theis, L.; Richardt, C.; and Yang, Y.-L. 2019. Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7588–7597.

Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.

Piccinelli, L.; Sakaridis, C.; and Yu, F. 2023. iDisc: Internal discretization for monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 21477–21487.

Po, R.; and Wetzstein, G. 2023. Compositional 3d scene generation using locally conditioned diffusion. *arXiv preprint arXiv:2303.12218*.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.

Rockwell, C.; Fouhey, D. F.; and Johnson, J. 2021. Pixelsynth: Generating a 3d-consistent experience from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 14104–14113.

Roessle, B.; Barron, J. T.; Mildenhall, B.; Srinivasan, P. P.; and Nießner, M. 2022. Dense depth priors for neural radiance fields from sparse input views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12892–12901.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.

Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, 234–241. Springer.

Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115: 211–252.

Sargent, K.; Koh, J. Y.; Zhang, H.; Chang, H.; Herrmann, C.; Srinivasan, P.; Wu, J.; and Sun, D. 2023a. Vq3d: Learning a 3d-aware generative model on imagenet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 4240–4250.

Sargent, K.; Li, Z.; Shah, T.; Herrmann, C.; Yu, H.-X.; Zhang, Y.; Chan, E. R.; Lagun, D.; Fei-Fei, L.; Sun, D.; et al. 2023b. Zeronvs: Zero-shot 360-degree view synthesis from a single real image. *arXiv preprint arXiv:2310.17994*.

Schönberger, J. L.; and Frahm, J.-M. 2016. Structure-from-Motion Revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.

Schönberger, J. L.; Zheng, E.; Pollefeys, M.; and Frahm, J.-M. 2016. Pixelwise View Selection for Unstructured Multi-View Stereo. In *European Conference on Computer Vision (ECCV)*.

Schuhmann, C.; Beaumont, R.; Vencu, R.; Gordon, C.; Wightman, R.; Cherti, M.; Coombes, T.; Katta, A.; Mullis, C.; Wortsman, M.; et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35: 25278–25294.

Song, L.; Chen, A.; Li, Z.; Chen, Z.; Chen, L.; Yuan, J.; Xu, Y.; and Geiger, A. 2023. Nerfplayer: A streamable dynamic scene representation with decomposed neural radiance fields. *IEEE Transactions on Visualization and Computer Graphics*, 29(5): 2732–2742.

Song, Y.; Sohl-Dickstein, J.; Kingma, D. P.; Kumar, A.; Ermon, S.; and Poole, B. 2020. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*.

Van Den Oord, A.; Vinyals, O.; et al. 2017. Neural discrete representation learning. *Advances in neural information processing systems*, 30.

Voleti, V.; Yao, C.-H.; Boss, M.; Letts, A.; Pankratz, D.; Tochilkin, D.; Laforte, C.; Rombach, R.; and Jampani, V. 2024. Sv3d: Novel multi-view synthesis and 3d generation from a single image using latent video diffusion. *arXiv preprint arXiv:2403.12008*.

Wang, G.; Wang, P.; Chen, Z.; Wang, W.; Loy, C. C.; and Liu, Z. 2023a. PERF: Panoramic Neural Radiance Field from a Single Panorama. *arXiv preprint arXiv:2310.16831*.

Wang, S.; Saharia, C.; Montgomery, C.; Pont-Tuset, J.; Noy, S.; Pellegrini, S.; Onoe, Y.; Laszlo, S.; Fleet, D. J.; Soricut, R.; et al. 2023b. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 18359–18369.

Wang, Z.; Lu, C.; Wang, Y.; Bao, F.; Li, C.; Su, H.; and Zhu, J. 2024. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *Advances in Neural Information Processing Systems*, 36.

Wu, J.; Zhang, C.; Xue, T.; Freeman, B.; and Tenenbaum, J. 2016. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *Advances in neural information processing systems*, 29.

Yang, B.; Luo, Y.; Chen, Z.; Wang, G.; Liang, X.; and Lin, L. 2023a. Law-diffusion: Complex scene generation by diffusion with layouts. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 22669–22679.

Yang, G.; Huang, X.; Hao, Z.; Liu, M.-Y.; Belongie, S.; and Hariharan, B. 2019. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF international conference on computer vision*, 4541–4550.

Yang, Y.; Liu, W.; Yin, F.; Chen, X.; Yu, G.; Fan, J.; and Chen, T. 2023b. VQ-NeRF: Vector Quantization Enhances Implicit Neural Representations. *arXiv preprint arXiv:2310.14487*.

Yang, Y.; Yin, F.; Liu, W.; Fan, J.; Chen, X.; Yu, G.; and Chen, T. 2024. PM-INR: Prior-Rich Multi-Modal Implicit Large-Scale Scene Neural Representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 6594–6602.

Yin, F.; Huang, Z.; Chen, T.; Luo, G.; Yu, G.; and Fu, B. 2023. Dcnet: Large-scale point cloud semantic segmentation with discriminative and efficient feature aggregation. *IEEE Transactions on Circuits and Systems for Video Technology*.

Yin, F.; Liu, W.; Huang, Z.; Cheng, P.; Chen, T.; and YU, G. 2022. Coordinates Are NOT Lonely–Codebook Prior Helps Implicit Neural 3D Representations. *arXiv preprint arXiv:2210.11170*.

Yin, F.; and Zhou, S. 2020. Accurate estimation of body height from a single depth image via a four-stage developing network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8267–8276.

Yin, W.; Zhang, J.; Wang, O.; Niklaus, S.; Mai, L.; Chen, S.; and Shen, C. 2021. Learning to recover 3d scene shape from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 204–213.

Yu, A.; Ye, V.; Tancik, M.; and Kanazawa, A. 2021. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4578–4587.

Yu, H.-X.; Duan, H.; Hur, J.; Sargent, K.; Rubinstein, M.; Freeman, W. T.; Cole, F.; Sun, D.; Snavely, N.; Wu, J.; et al. 2023. WonderJourney: Going from Anywhere to Everywhere. *arXiv preprint arXiv:2312.03884*.

Zhang, F.; Zhang, Y.; Zheng, Q.; Ma, R.; Hua, W.; Bao, H.; Xu, W.; and Zou, C. 2024a. 3D-SceneDreamer: Text-Driven 3D-Consistent Scene Generation. *arXiv preprint arXiv:2403.09439*.

Zhang, J.; Li, X.; Wan, Z.; Wang, C.; and Liao, J. 2024b. Text2nerf: Text-driven 3d scene generation with neural radiance fields. *IEEE Transactions on Visualization and Computer Graphics*.

Zhang, S.; Han, Z.; Lai, Y.-K.; Zwicker, M.; and Zhang, H. 2019. Stylistic scene enhancement GAN: mixed stylistic enhancement generation for 3D indoor scenes. *The Visual Computer*, 35: 1157–1169.

# A    Appendix

## A.1   Implementation Details

We implement our Scene123 with the Pytorch framework (Paszke et al. 2019) and adopt TensoRF (Chen et al. 2022) as the radiance field. To ensure TensoRF can accommodate scene generation over a large view range, we position the camera near the center of the NeRF bounding box and configure it with outward-facing viewpoints. The dimension of the masked codebook is $2048 \times 16$. For only text prompt input, we utilize the stable diffusion model in version 2.0 (Rombach et al. 2022) to generated initial image $\mathbf{I_0}$ related to the input prompt. Moreover, for depth estimation, we use the boosting monocular depth estimation method (Miangoleh et al. 2021) with pre-trained LeReS model (Yin et al. 2021) to estimate the depth for each view. For the image-to-video pipeline, we opt for the Stable Video Diffusion in version SVD-XT, which is the same architecture as SVD (Blattmann et al. 2023a) but finetuned for 25 frame generation. During training, we use the same setting as (Chen et al. 2022) for the optimizer and learning rate and set the hyperparameters in our objective function as $\lambda_d = 0.005$, $\lambda_t = 0.001$, $\lambda_{dist} = 0.001$.

**GAN-based Training Strategy Details.** For the incorporated discriminator, we adopt a similar architecture, regularization function, and loss weight as the EG3D model (Chan et al. 2022), with some distinctions. The vanilla structure of a 3D GAN involves a 3D generator that incorporates a super-resolution module, accompanied by a discriminator that accepts both coarse and fine image inputs (Chan et al. 2022). Notably, due to the contextual disparities between text-to-3D and 3D GAN applications, we opt to omit the super-resolution component from the 3D GAN architecture. This choice stems from the consistent need to extract mesh or voxel representations from the 3D model within the context of text-to-3D.

The camera motion in the SVD model is sometimes limited (Blattmann et al. 2023a), making it unsuitable for directly training the NeRF model. Instead, we use our 3D-Aware Generative Refinement approach, which leverages a discriminator to optimize the process. In this setup, the video diffusion model's output is treated as real data, while the 3D neural radiance field model's renderings are treated as fake data. Despite the small range of camera motion in the images generated by SVD, the adversarial training process provides necessary regularization. The discriminator learns to distinguish between real support set images and fake renderings, thereby pushing the neural radiance field model to produce outputs that are indistinguishable from real images. This continuous adversarial interaction compels the generator to improve by learning finer details and more accurate textures, ultimately enhancing the quality of the 3D scene generation. This method ensures that even with limited camera motion, the GAN loss remains effective and contributes to generating high-quality 3D scenes, as shown in Fig. 7 in the manuscript.

**Depth alignment details.** We use the depth estimation model $f_e$ to estimate the depth map $D_k^E$ for the initial view $\mathbf{I_0}$. Note that, unlike the depth map $D_0$ of the initial view, $D_k^E$ cannot be directly taken as the supervision to update the radiance field since it is predicted independently and could conflict with known depth maps such as $D_k^R$ in the overlapping regions. To solve this issue, we implement depth alignment to align the estimated depth map to the known depth values in the radiance field (Zhang et al. 2024b).

Due to the lack of geometric constraint during the depth estimation, the predicted depth values could be misaligned in the overlapping regions (Luo et al. 2020), for example, the estimated depth $D_k^E$ of the inpainted view may be inconsistent with the depth $D_k^R$ rendered from NeRF since $D_k^R$ is constrained by previous known views. The inconsistency is manifested in two aspects: scale difference and value difference. For instance, the *distance difference* of

two pixel-aligned spatial points and the *depth value* of a specific point could be both different in depth maps estimated from different views. The former is the scale difference and the latter is the value difference. In the case of scale difference, we cannot align both points by shift processing because even if we align the depth value of one of the points, the other point is still misaligned. To eliminate the scale and value differences between the overlapping regions of the rendered depth map $D_k^R$ and the estimated depth map $D_k^E$ of the novel view, we introduce a two-stage depth alignment strategy. Specifically, we first globally align these two depth maps by compensating for mean scale and value differences. Then we finetune a pre-trained depth alignment network to produce a locally aligned depth map.

To determine the mean scale and value differences, we first randomly select $M$ pixel pairs from the overlapping regions and deduce their 3D positions under depth $D_k^R$ and $D_k^E$, denoted as $\left\{ (\mathbf{x}_j^R, \mathbf{x}_j^E) \right\}_{j=1}^M$. Next, we calculate the average scaling score $s$ and depth offset $\delta$ to approximate the mean scale and value differences:

$$s = \frac{1}{M-1} \sum_{j=1}^{M-1} \frac{\| \mathbf{x}_j^R - \mathbf{x}_{j+1}^R \|_2}{\| \mathbf{x}_j^E - \mathbf{x}_{j+1}^E \|_2}, \tag{6}$$

$$\delta = \frac{1}{M} \sum_{j=1}^{M} \left( z\left( \mathbf{x}_j^R \right) - z\left( \hat{\mathbf{x}}_j^E \right) \right), \tag{7}$$

where $\hat{\mathbf{x}}_j^E = s \cdot \mathbf{x}_j^E$ indicates the scaled point and $z(\mathbf{x})$ represents the depth value of point $\mathbf{x}$. Then $D_k^E$ can be globally aligned with $D_k^R$ by $D_k^{global} = s \cdot D_k^E + \delta$.

Since depth maps used in our pipeline are predicted by a network, the differences between $D_k^R$ and $D_k^E$ are not linear, that is why the global depth aligning process cannot solve the misalignment problem. To further mitigate the local difference between $D_k^{global}$ and $D_k^R$, we train a pixel-to-pixel network $f_\psi$ for nonlinear depth alignment. During optimization of each view, we optimize the parameter $\psi$ of the pre-trained depth alignment network $f_\psi$ by minimizing their least square error in the overlapping regions:

$$\min_{\psi} \left\| \left( f_\psi(D_k^{global}) - D_k^R \right) \odot M_k \right\|_2. \tag{8}$$

Finally, we can derive the locally aligned depth using the optimized depth alignment network: $\hat{D}_k = f_\psi(D_k^{global})$.

## A.2   Additional analysis

**Comparison between SD inpainting and Our Consistency-Enhanced MAE in detail completion.** Our Consistency-Enhanced MAE ensures consistency across views by using a learnable codebook to distribute global information to every invisible area. Traditional VAE methods quantize the latent space, leading to artifacts at the boundaries between quantized latents. In contrast, the cross attention mechanism allows for dynamic and flexible feature weighting, better captures global context, and preserves detailed information without the loss associated with quantization. The learnable codebook continuously adapts during training, enhancing the model's stability and performance in generating high-quality 3D scenes.

However, the stable diffusion (SD) model, which is originally utilized to inpaint images in LucidDreamer (Chung et al. 2023), is less effective for the multi-view complementation task. While SD can achieve complementation through conditional diffusion, it struggles with complex structures and view-consistency tasks, as it is not tailored for this purpose. Although SD performs well in generating new images, it may not be as effective in complementing local details compared to MAE, which is specifically designed for
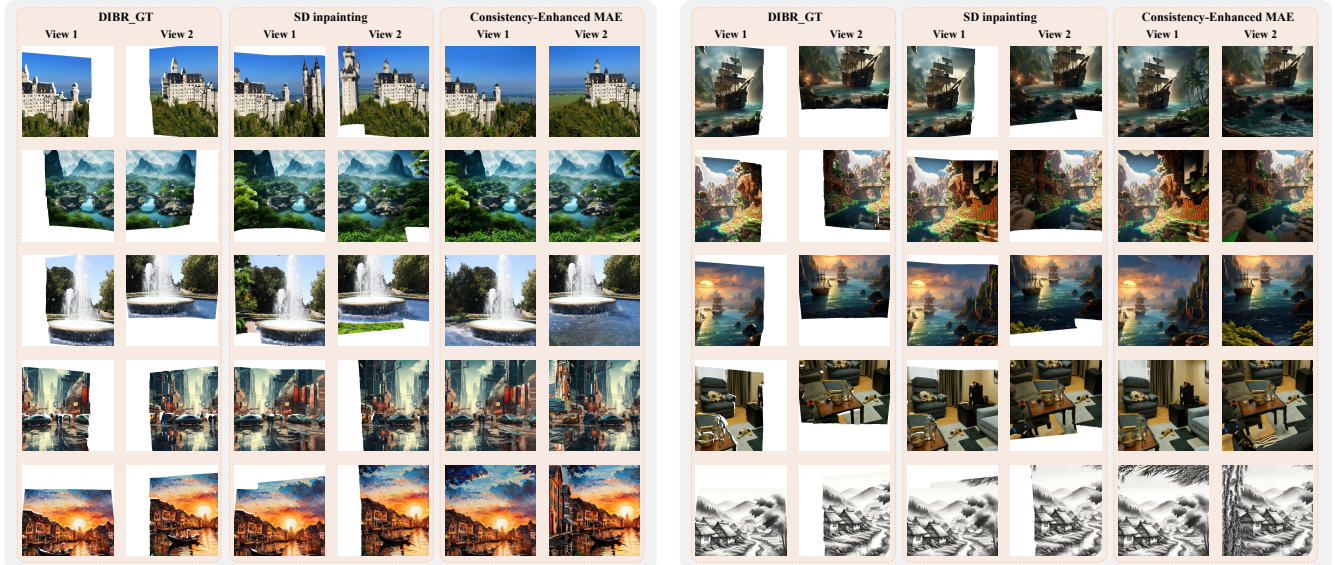
Figure 8: Comparisons of SD inpainting and our Consistency-Enhanced MAE in handling detailed complementation.

this task. This superiority of our Consistency-Enhanced MAE in handling detailed complementation is evident in Fig. 8.

**Video Generation Model promotes our Scene123's quality.** The 3D scenes generated through our 3D-Aware Generative Refinement module are of higher quality than the image-to-video generation pipeline. As depicted in Fig. 9, our Scene123 does not strictly require a high-quality support set. The novel views generated by our method are generally superior to those produced by the image-to-video pipeline. This is because our 3D-Aware Generative Refinement uses a discriminator to manage optimization, distinguishing between real and fake images. This adversarial interaction improves the generator's ability to learn finer details and more accurate textures, enhancing the overall quality of 3D scene generation. This demonstrates our approach's adaptability in handling image-to-video generation failures.

## A.3 User Study

For completeness, we follow previous works (Wang et al. 2024; Yu et al. 2023) and conduct a user study by comparing Scene123 with WonderJourney (Yu et al. 2023), LucidDreamer (Chung et al. 2023) under 40 image prompts, 20 image prompts for each baseline. For text prompt input, we also compare Scene123 with Wonder-Journey (Yu et al. 2023), LucidDreamer (Chung et al. 2023) and Text2NeRF (Zhang et al. 2024b) under 60 text prompts, 20 text prompts for each baseline. The participants are shown the generated results of our Scene123 and baselines and asked to choose the better one in terms of fidelity, details and vividness. We collect results from 39 participants, yielding 3120 pairwise comparisons. The results are shown in Tab. 4. Our method outperforms all of the baselines.

| Method | WonderJourney | LucidDreamer | Text2NeRF |
|---|---|---|---|
| Prefer baseline | 31.12 | 23.42 | 46.27 |
| Prefer Scene123 (Ours) | **68.88** | **76.58** | **53.73** |

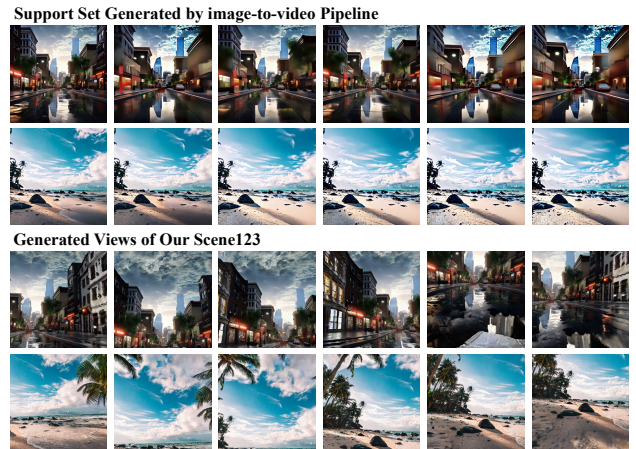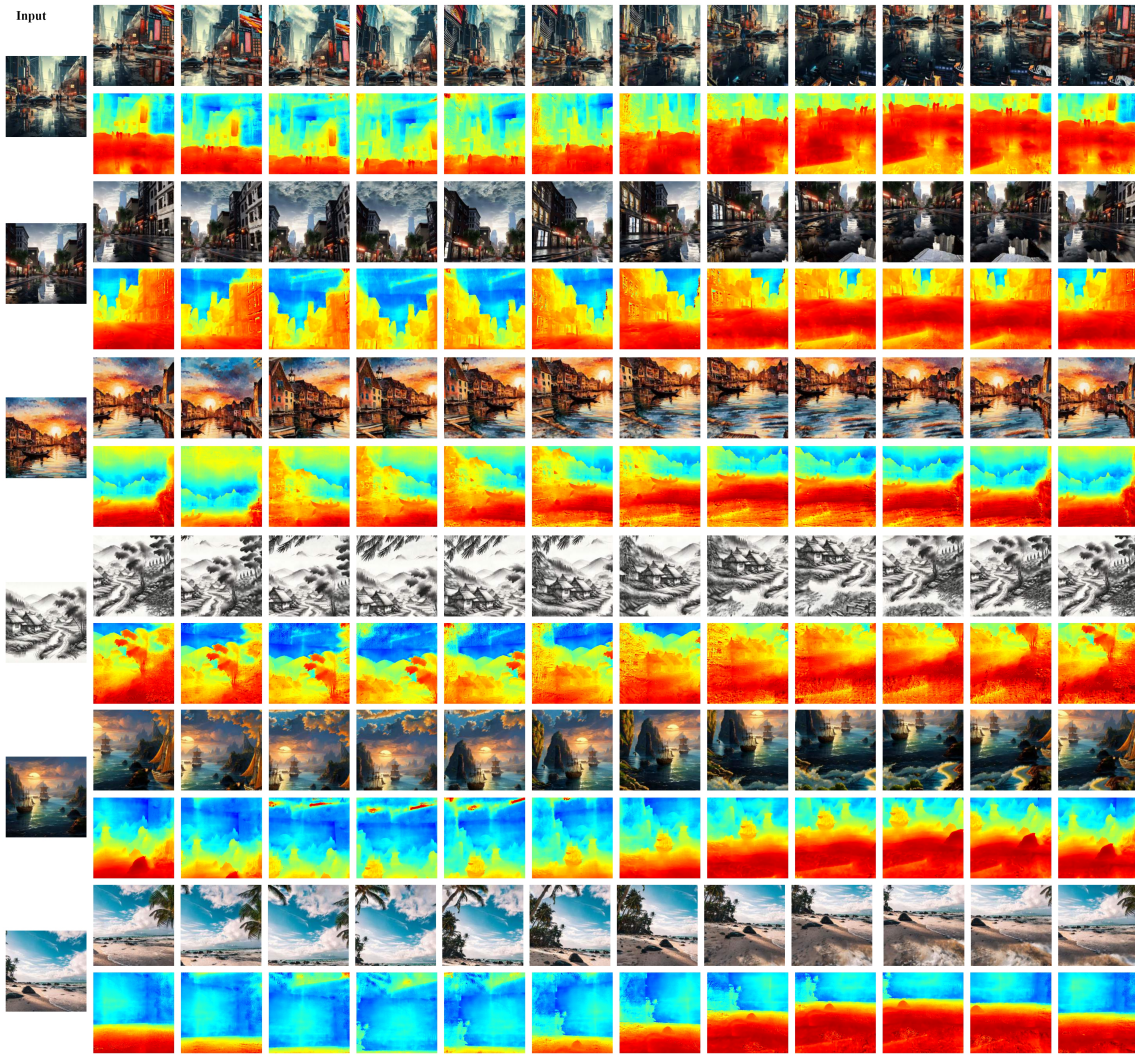Table 4: Results of user study. The percentage of user preference (↑) is reported in the table.



Figure 9: Comparisons between support set generated by image-to-video pipeline and the generated views of Our Scene123.
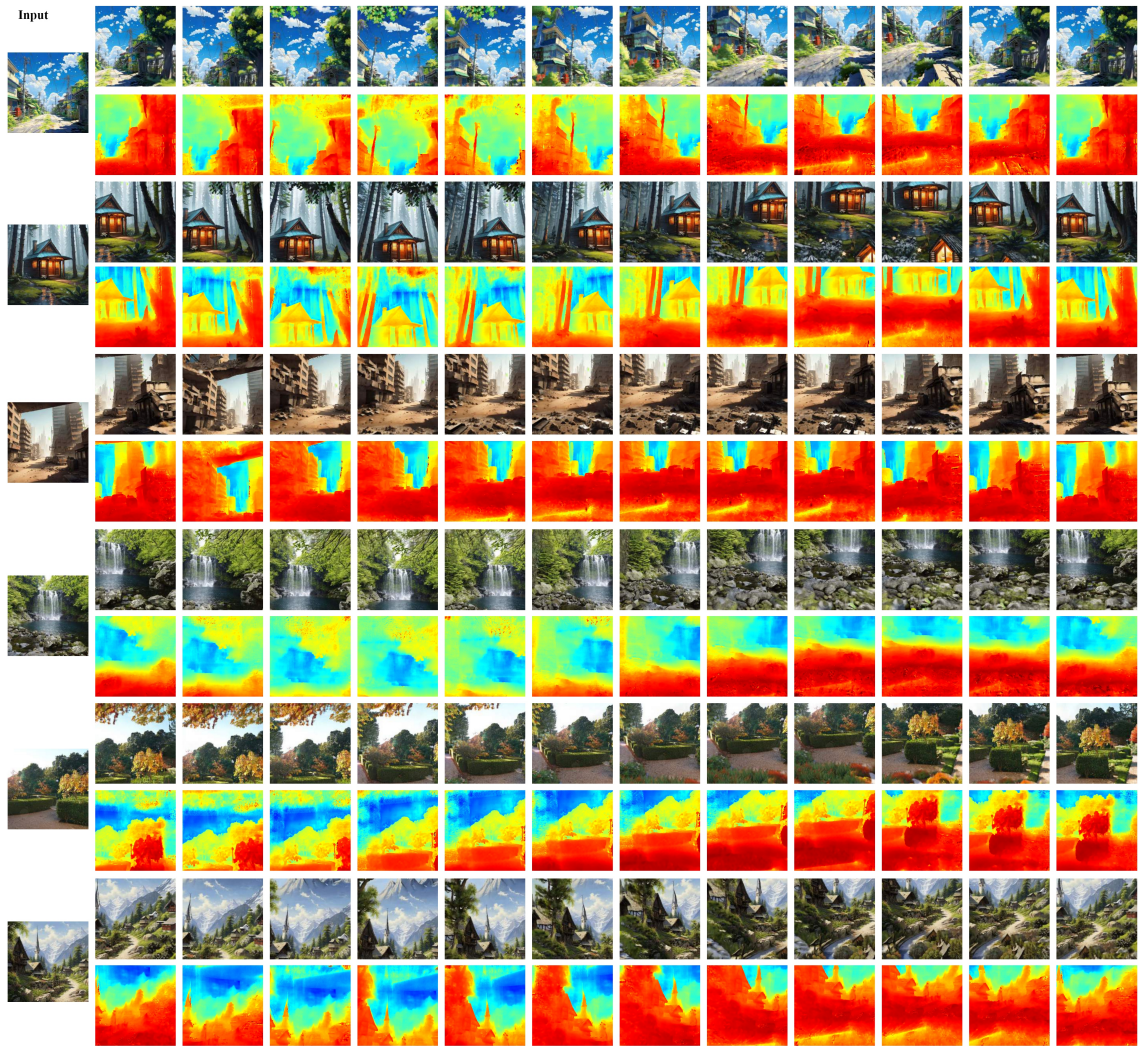
## A.4 More Qualitative Results

We provide more qualitative results in Fig. 10, Fig. 11, Fig. 12, Fig. 13, including indoor scenes, outdoor scenes, outdoor buildings, object-centered scenes with realistic renderings and precise depth details.
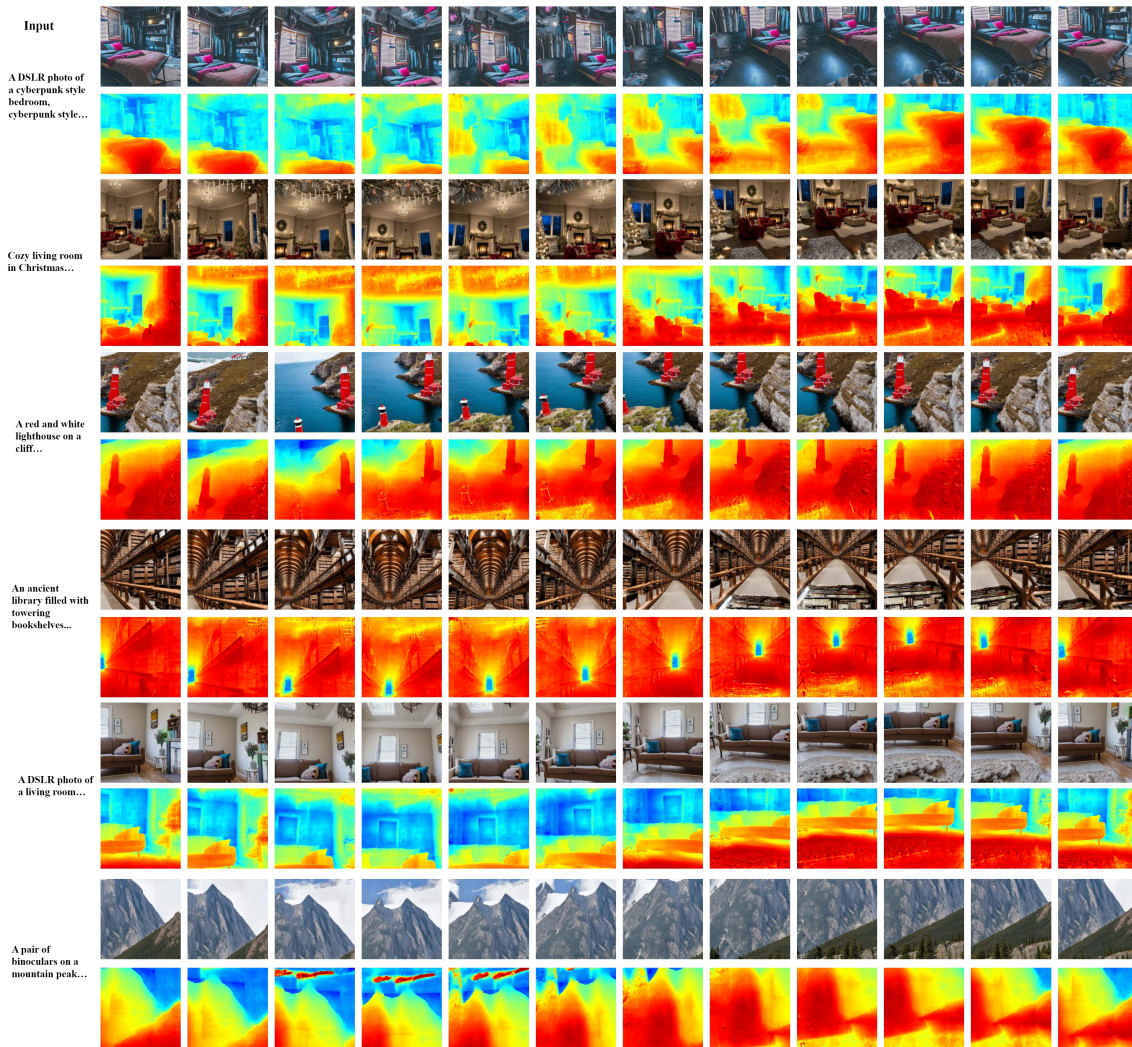
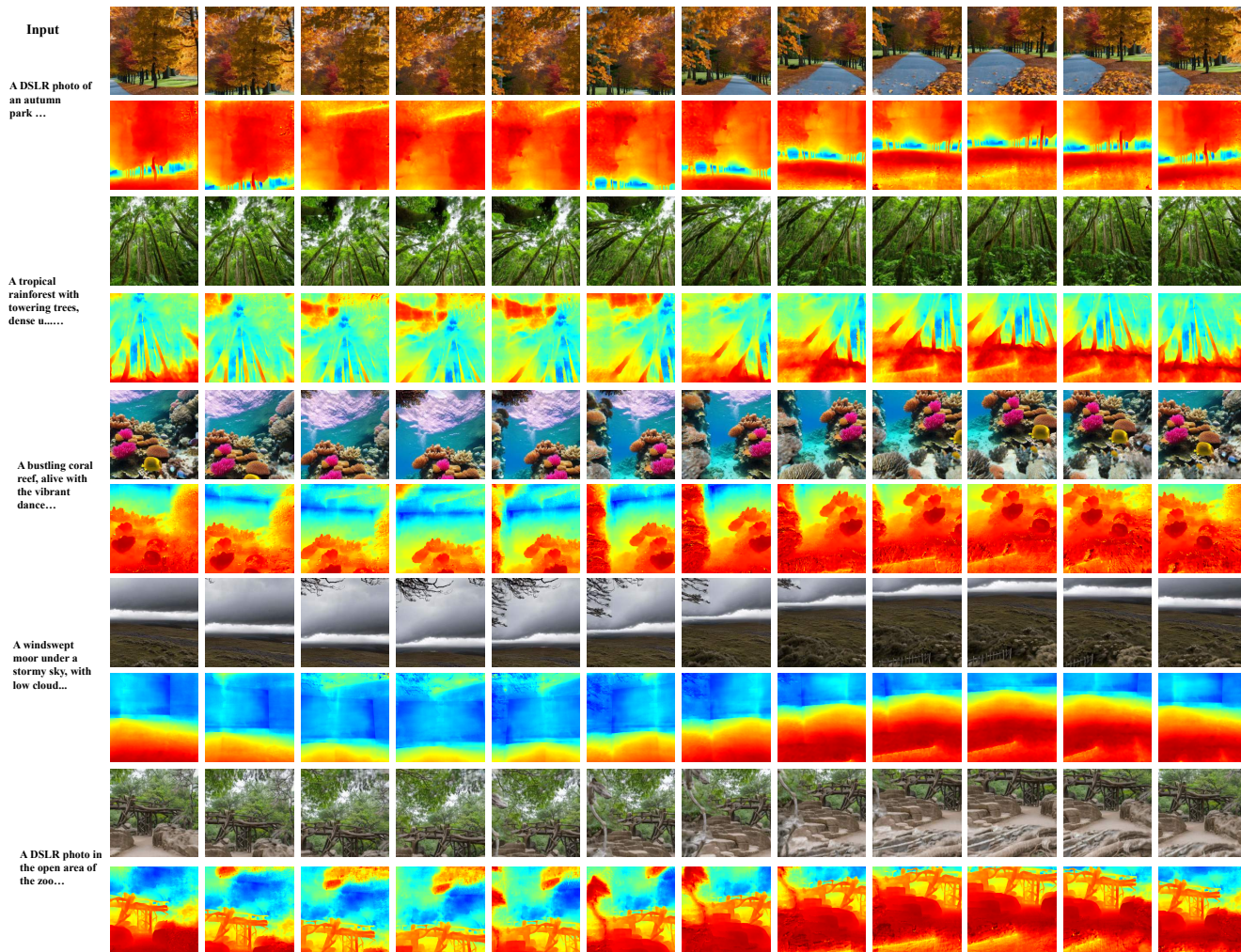Figure 10: More qualitative examples generated by our Scene123 from a single image input.

Part 2 / 4

Figure 11: More qualitative examples generated by our Scene123 from a single image input.

Figure 12: More qualitative examples generated by our Scene123 from a text prompt input.

Part 4 / 4

Figure 13: More qualitative examples generated by our Scene123 from a text prompt input.