

InFusion: Inpainting 3D Gaussians via Learning Depth Completion from Diffusion Prior

Zhiheng Liu^{1,4} *, Hao Ouyang^{2,3} *, Qiuyu Wang³, Ka Leong Cheng^{2,3}, Jie Xiao^{1,4}, Kai Zhu⁴, Nan Xue³, Yu Liu¹, Yujun Shen³, and Yang Cao¹ †

¹ University of Science and Technology of China

² The Hong Kong University of Science and Technology

³ Ant Group

⁴ Alibaba Group

Abstract. 3D Gaussians have recently emerged as an efficient representation for novel view synthesis. This work studies its editability with a particular focus on the inpainting task, which aims to supplement an incomplete set of 3D Gaussians with additional points for visually harmonious rendering. Compared to 2D inpainting, the crux of inpainting 3D Gaussians is to figure out the rendering-relevant properties of the introduced points, whose optimization largely benefits from their initial 3D positions. To this end, we propose to guide the point initialization with an image-conditioned depth completion model, which learns to directly restore the depth map based on the observed image. Such a design allows our model to fill in depth values at an aligned scale with the original depth, and also to harness strong generalizability from large-scale diffusion prior. Thanks to the more accurate depth completion, our approach, dubbed **InFusion**, surpasses existing alternatives with sufficiently better fidelity and efficiency (*i.e.*, $\sim 20\times$ faster) under various complex scenarios. We further demonstrate the effectiveness of **InFusion** with several practical applications, such as inpainting with user-specific texture or with novel object insertion. Our code is public available at <https://johanan528.github.io/Infusion/>.

Keywords: Gaussian splatting · 3D inpainting · Monocular depth completion

1 Introduction

Recent developments in 3D representation [4, 45, 65, 96] have highlighted 3D Gaussians [14, 45, 93, 102, 107] as an essential approach for novel view synthesis, owing to the ability to produce photorealistic images with impressive rendering speed. 3D Gaussians offer explicit representation and the capability for real-time processing, which significantly enhances the practicality of editing 3D scenes. The study of how to editing 3D Gaussians is becoming increasingly vital, particularly for interactive downstream applications such as virtual and augmented reality

*These authors contributed equally to this work. †Corresponding author.

(VR/AR). Our research focuses on the inpainting tasks that are crucial for the seamless integration of edited elements, effectively filling in missing parts and serving as a foundational operation for further manipulations.

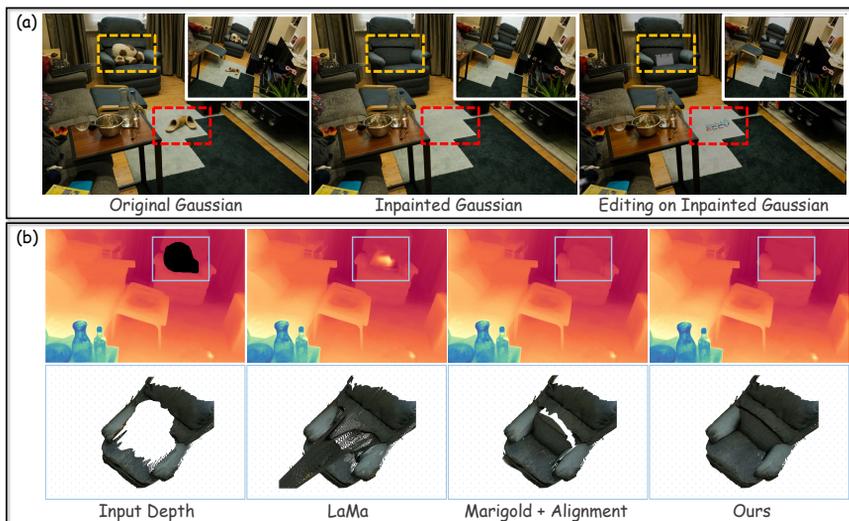


Fig. 1: We present **InFusion**, an innovative approach that delivers efficient, photorealistic inpainting for 3D scenes with 3D Gaussians. As demonstrated in (a), **InFusion** enables the seamless removal of 3D objects, along with user-friendly texture editing and object insertion. Illustrated in (b), **InFusion** learns depth completion with diffusion prior, significantly enhancing the depth inpainting quality for general objects. We show the visualizations of the unprojected points, which exhibit substantial improvements over baseline models [44, 92].

Initial explorations into 3D Gaussian inpainting have focused on growing Gaussians from the boundary of the unpainted regions, using inpainted 2D multiview images for guidance [13, 29, 106]. This method, however, tends to produce blurred textures due to inconsistencies in the generation process, and the growing can be quite slow. Notably, the training quality for Gaussian models is significantly improved when the initial points are precisely positioned within the 3D scene, particularly on object surfaces. A practical solution to improve the fine-tuning of Gaussians is to predetermine these initial points where inpainting will occur, thereby simplifying the overall training process. In allocating initial points for Gaussian inpainting, the role of depth map inpainting can be pivotal. The ability to convert inpainted depth maps into point clouds facilitates a seamless transition to 3D space, while also leveraging the potential to train on expansive datasets [62, 63, 84].

To this end, we introduce **InFusion**, an innovative approach to 3D Gaussian inpainting that leverages depth completion learned from diffusion models [1, 9, 75, 79]. Our method demonstrates that with a robustly learned depth inpainting model, we can accurately determine the placement of initial points, significantly

elevating both the fidelity and efficiency of 3D Gaussian inpainting. In particular, we first inpaint the depth in the reference view, then unproject the points into the 3D space to achieve optimal initialization. However, current depth inpainting methodologies [44, 67, 92, 106] are often a limiting factor; commonly, they lack the generality required to accurately complete object depth, or they produce depth maps that misalign with the original, with errors amplified during unprojection. In this work, we harness the power of pre-trained latent diffusion models, training our depth inpainting model with diffusion-based priors to substantially enhance the quality of our inpainting results. The model exhibits a marked improvement in aligning with the unpainted regions and in reconstructing the depth of objects. This enhanced alignment capability ensures a more coherent extension of the existing geometry into the inpainted areas, leading to a seamless integration within the 3D scene. Furthermore, to address challenging scenarios involving large occlusions, we design InFusion with a progressive strategy that showcases its capability to resolve such complex cases.

Our extensive experiments on various datasets, which include both forward-facing and unbounded 360-degree scenes, demonstrate that our method outperforms the baseline approaches in terms of visual quality and inpainting speed, being 20 times faster. With the effective depth inpainting framework based on a pre-trained LDM, we demonstrate that the integration of 3D Gaussians with depth inpainting offers an efficient and feasible approach to completing 3D scenes. The strength of LDMs [75, 79] is pivotal to our approach, allowing our model to inpaint not just the background but also to complete objects. Beyond the core functionality, our method facilitates additional applications, such as user-interactive texture inpainting, which enhances user engagement by allowing direct input into the inpainting process. We also demonstrate the adaptability of our method for downstream tasks, including scene manipulation and object insertion, revealing the broad potential of our approach in the context of editing and augmenting 3D spaces.

2 Related Work

2.1 Image and Video Inpainting

Image and video inpainting is an important editing task [18, 48, 70, 72, 75, 81, 104, 109] that aims to restore the missing regions within an image or video by inferring visually consistent content. Traditional works for image inpainting [2, 3, 5, 23, 27, 94] typically involve extracting low-level features to restore damaged areas. Similarly, in video inpainting [33, 38, 69, 70, 86, 91, 101], the restoration process is often approached as an optimization task based on patch sampling. However, these methods generally lack capacity when handling images with large missing regions or corrupted videos with complex motions. Recently, deep learning has not only empowered inpainting models to overcome these challenges in restoration but has also expanded their capacity to generate new, semantically plausible content [76]. State-of-the-art image inpainting methods [18, 24, 48, 52, 60, 75, 81, 92, 108] excel at effectively handling large mask inpainting tasks on high-resolution images;

cutting-edge techniques on video inpainting [5, 31, 51, 53, 56, 104, 111, 112, 117, 118] commonly leverage flow-guided propagation and video Transformers to restore missing parts in videos with natural and spatiotemporally coherent content.

2.2 3D Scene Inpainting

With the increasing accessibility of 3D reconstruction models, there is a growing demand for 3D scene editing [16, 35, 49, 73, 110, 114, 120]. 3D scene inpainting is one prominent application to fill in the missing parts within a 3D space, such as removing objects from the scene and generating plausible geometry and texture to complete the inpainted regions. Early inpainting works mainly focuses on performing geometry completion [20–22, 34, 42, 43, 74, 90, 97, 103]. Recent advancements in 3D inpainting techniques have facilitated the simultaneous inpainting of both semantics and geometry by successfully handling the interplay between these two aspects [100]. They can be broadly categorized into two groups based on the adopted 3D representation: NeRF [64] and Gaussian Splatting (GS) [46]. Some NeRF-based methods [47, 49, 58, 68, 87] leverage CLIP [77] or DINO features [10] to learn 3D semantics for inpainting; others [15, 55, 66, 67, 95, 98, 99] typically rely on 2D image inpainting models with depth or segmentation priors to optimize NeRFs through neural fields rendering. In contrast to inpainting on NeRF, several methods [13, 39, 41, 106] explore inpainting techniques on GS models, thanks to their notable advantages such as impressive rendering efficiency and high-quality reconstruction. In our paper, we further improve the efficiency and the quality of 3D inpainting within GS settings.

2.3 Diffusion Models for Monocular Depth

The explicit nature of 3D Gaussians makes the accurate allocation of inpainted points within 3D scenes (*e.g.*, object surfaces) highly beneficial for 3D scene inpainting via optimization. A direct and effective solution is to utilize the 2D depth prior of reference views obtained through monocular depth estimation [8, 28, 32, 54, 78, 105, 119] or completion [6, 30, 57, 59, 61, 85, 113, 115] models to initialize the inpainted 3D points. Thanks to the superior performance of latent diffusion models (LDM) [7, 9, 36, 75, 79, 80, 88, 89], it opens up the possibility of enhancing depth learning by leveraging or distilling the capabilities of these models. Several methods have attempted to employ diffusion priors for estimating monocular depth [26, 40, 44, 82, 83, 116]. However, learning from LDM for monocular depth completion (or inpainting) receives less attention. While some methods [67, 99] employ LaMa [92] to inpaint depth in the Jet color space, the precision of the resulting inpainted depth map is compromised due to the lossy quantization process when converting metric depth to the Jet color space. To the best of our knowledge, our work is the first resolve this problem by training an accurate depth completion model from diffusion prior [75].

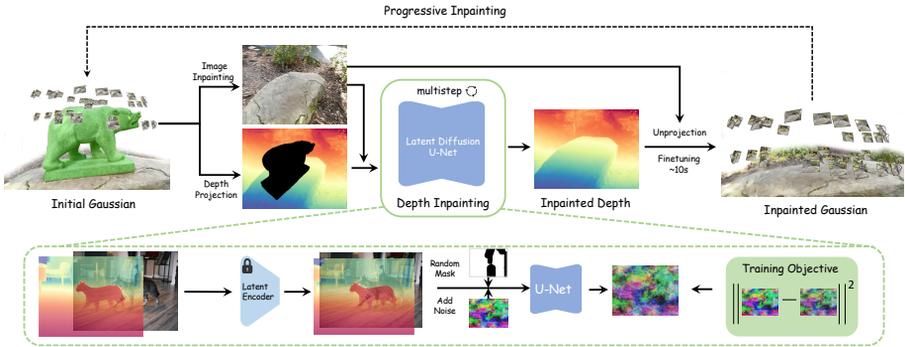


Fig. 2: Illustration of Infusion driven by Depth Inpainting. *Top:* To remove a target from the optimized 3D Gaussians, our InFusion first inpaints a selected one-view RGB image and applies the proposed diffusion model for depth inpainting to the depth projection of the targeted 3D Gaussians. The progressive scheme addresses view-dependent occlusion issues by utilizing other unobstructed viewpoints. *Bottom:* A detailed view of the training pipeline for the depth inpainting U-Net is presented. We employ a mask-driven denoising diffusion for training of the U-Net, which utilizes a frozen latent tokenizer by taking the RGB image and depth map as inputs.

3 Method

3.1 Overview

Formally, 3D scenes can be represented by 3D Gaussians Θ , given a collection of multi-view images $\mathcal{I} = \{I_i\}_{i=1}^n$, accompanied by respective camera poses $\Pi = \{\pi_i\}_{i=1}^n$ [46]. Our objective is to edit the scene Θ with a particular focus on inpainting, which aims to supplement an incomplete set of 3D Gaussians. The complexity of 3D Gaussian inpainting arises due to potential inconsistencies in the supervision provided by the 2D inpainted images from multiview. Nevertheless, three key observations inspire our solution design to address the challenges:

- The reconstruction quality of the optimized 3D Gaussians for novel view synthesis is highly sensitive to the initialization, especially when the view number is limited. Hence, we are motivated to carefully place the initial points within the inpainting regions for enhancing the inpainting quality.
- Contemporary research indicates that the initialization of 3D Gaussians with unprojected depth maps [12, 19, 71] yields promising results due to explicitness. This observation implies that using inpainted depth images for initializing the missing region could be advantageous.
- Incorporating a diffusion prior [1, 9, 44, 75, 79] into depth estimation markedly improved accuracy especially for general objects. This finding indicates that a similar approach can be adopted to leverage diffusion priors for benefiting depth inpainting.

Leveraging the key observations discussed earlier, our pipeline is illustrated in Figure 2. Starting with the 3D Gaussians Θ , we first segment out and discard unwanted Gaussians under the guidance of masks $\mathcal{M} = \{m_i\}_{i=1}^n$, which delineate the targeted regions for modification. As mentioned, depth inpainting can play a crucial role in determining the initial placement of Gaussians. To achieve this, we select a reference view and perform inpainting on both the image and its corresponding depth map to facilitate accurate unprojection. Existing depth inpainting models may not possess the versatility needed for precise depth completion or may produce depths that are inconsistent with the unpainted regions. Such misalignments lead to suboptimal inpainting outcomes. To address this, we develop a more generalized depth inpainting model that harnesses the strengths of natural diffusion processes. In situations with substantial occlusion, relying on a single reference view may prove insufficient. To solve this, our approach incorporates multiple reference views through a progressive inpainting strategy.

The remainder of our methods are structured as the following. We describe the specifics of the diffusion-based depth completion model in Sec. 3.2 and use this model to do 3D scene inpainting in Sec. 3.3. Finally, we provide the details of progressive inpainting in Sec. 3.4.

3.2 Diffusion Models for Depth Completion

A precise and reliable depth inpainting model is essential to obtain a well-founded set of initial points for inpainting Gaussians. We build our depth completion model on latent diffusion models (LDMs) [79] for the strong priors due to their training on extensive, internet-scale collections of images. Given a set of color images and their corresponding depth, as well as various random masks, we seek to learn a model with the ability to inpaint the masked depth. The following three sections describe our diffusion-based depth completion model in details.

Diffusion Models We formulate depth completion as a task of conditional denoising diffusion generation. The LDMs operates by conducting diffusion processes within a lower-dimensional latent space, facilitated by a pre-trained Variational Auto-Encoder (VAE) \mathcal{E} . Diffusion steps are performed on these noisy latents where a denoising U-Net ϵ_θ iteratively removes noise to get clean latents. During inference, the U-Net is applied to denoise pure Gaussian noise into a clean latent. The image recovery is then achieved by passing these refined latents through the VAE decoder \mathcal{D} . This ensures that the depth completion model benefits from the powerful generative capabilities inherent in LDMs while also maintaining efficiency by operating within a compressed latent space.

Training We develop our model on top of a pre-trained text-to-image LDM (Stable Diffusion [79]) to save computational resources and enhance training efficiency. Modifying the existing model architecture, we adapt it for image-conditioned depth completion tasks. An outline of the refined fine-tuning process is presented in Fig. 2.

Our depth completion diffusion model accepts a trio of inputs: a depth map d , a corresponding color image I , and a mask m . Leveraging the frozen VAE,

we encode both the color image and the depth map into a latent space, which serves as the foundation for training our conditional denoiser. To accommodate the VAE encoder’s design for 3-channel (RGB) inputs when presented with a single-channel depth map, we duplicate the depth information across three channels to create an RGB-like representation. We apply a linear normalization to ensure the depth values predominantly reside within the interval $[-1, 1]$ following Marigold [44], thereby conforming to the VAE’s expected input range. This normalization is executed via an affine transformation delineated as follows:

$$d' = \frac{d - d_2}{d_{98} - d_2} \times 2 - 1, \quad (1)$$

where d_2 and d_{98} represent the 2^{nd} and 98^{th} percentiles of individual depth maps, respectively. Such normalization facilitates the model’s concentration on affine-invariant depth completion, enhancing the robustness of the algorithm against scaling and translation.

The normalized depth d' and the color image are first encoded into the latent space with the encoder of the VAE:

$$z^{(d')} = \mathcal{E}(d'), z^{(I)} = \mathcal{E}(I), \quad (2)$$

The encoder produces a 4-channel feature map that has a lower resolution than the original input. To construct the image-conditioned depth completion model, we initially resize the mask m to align with the dimensions of $z^{(d')}$, yielding $m' = \text{downsample}(m)$. We then create a composite feature map by concatenating the noisy latent depth code $z_t^{(d')}$, the element-wise product of the clean latent depth code and the downscaled mask $z_m^{(d')} = z^{(d')} \odot m'$, and the latent image code $z^{(I)}$, along with m' , as follows:

$$z_t = \text{cat}(z_t^{(d')}, z_m^{(d')}, z^{(I)}, m'), \quad (3)$$

along the channel dimension, where $z_t^{(d')} = \alpha_t z^{(d')} + \sigma_t \epsilon$. The concatenated feature map z_t , comprising $4 + 4 + 4 + 1 = 13$ channels, is subsequently fed into the U-Net-based denoiser ϵ_θ .

At training time, U-Net parameters θ are updated by taking a data pair (I, d, m) from the training set, noising d with sampled noise ϵ at a random timestep t , computing the noise estimate $\hat{\epsilon} = \epsilon_\theta(z_t)$ and minimizing the denoising diffusion objective function following DDPM [36]:

$$\mathcal{L} = \mathbb{E}_{d, \epsilon, t} \|\epsilon - \epsilon_\theta(z_t)\|_2^2, \quad (4)$$

where $t \in \{1, 2, \dots, T\}$ indexes the diffusion timesteps, $\epsilon \in \mathcal{N}(0, I)$, and z_t , the noisy latent at timestep t , is calculated as Eq. (3).

Inference The inference of our depth completion model commences with an input comprising a depth map d , its corresponding color image I , and a mask m that delineates the target completion region. The color image I undergoes SDXL-based [75] image inpainting, resulting in $\tilde{I} = \mathcal{F}_I(I, m)$, where \mathcal{F}_I represents the

image inpainting model. Subsequently, we generate the concatenated feature map as defined in Eq. (3), which is then progressively refined according to the fine-tuning scheme. Leveraging the non-Markovian sampling strategy from DDIM [89] with re-spaced steps facilitates an accelerated inference. The final depth map is then derived from the latent representation decoded by the VAE decoder \mathcal{D} , followed by channel-wise averaging for post-processing.

3.3 Inpainting 3D Gaussians with Diffusion Priors

The trained diffusion model generates plausible depth completions, thereby serving as an effective initialization for the 3D Gaussians. Upon removing undesired points from 3D Gaussians, a set of reference views $\{I_{s(i_j)}\}_{j=1}^r$ is selected, where $s(i_j) \in \{1, 2, \dots, n\}$ and r denotes the total number of chosen views. For forward-facing and certain 360-degree inward-facing datasets, a single reference view ($r = 1$) is usually sufficient, whereas for more complex 360-degree scenes with occlusions, multiple reference views ($r > 1$) are required. In instances with $r > 1$, a progressive inpainting strategy is employed, detailed further in Sec. 3.4. The current discussion is focused on the $r = 1$ scenario.

Assuming without loss of generality, for $r = 1$, we designate the $s(i_1)^{th}$ view as the single reference. Initially, the color image $I_{s(i_1)} \odot m_{s(i_1)}$ is inpainted using an SDXL-based inpainting model to yield the restored image $\tilde{I}_{s(i_1)}$. The depth for the $s(i_1)^{th}$ view is then determined analogous to color rendering in GS:

$$d_{s(i_1)} = \sum_{i \in N_{s(i_1)}} z_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (5)$$

where z_i denotes the z-coordinate in the world coordinate system, and α_i represents the density of the corresponding point. It is important to note that the resulting depth d is incomplete, as it derives from Θ . To address this, we apply our diffusion-based depth completion model \mathcal{F}_d , which produces the refined depth map:

$$\tilde{d}_{s(i_1)} = \mathcal{F}_d(d_{s(i_1)}, \tilde{I}_{s(i_1)}, m_{s(i_1)}). \quad (6)$$

With the completed depth map $\tilde{d}_{s(i_1)}$, the inpainted image $\tilde{I}_{s(i_1)}$, and the corresponding camera pose $\Pi_{s(i_1)}$, we unproject $\tilde{d}_{s(i_1)}$ and $\tilde{I}_{s(i_1)}$ from image space to 3D coordinates to form a colored point cloud $\mathcal{P}_{s(i_1)}$. This point cloud is then merged with the original 3D Gaussian point cloud to achieve a robust initialization Θ' for subsequent GS fine-tuning.

Ultimately, the preliminary Gaussian model Θ' are fine-tuned merely 50~150 iterations to yield the final Gaussian model $\tilde{\Theta}$, using solely the selected view image $I_{s(i_1)}$. The optimization is also guided by \mathcal{L}_1 combined with D-SSIM at the $s(i_1)^{th}$ view:

$$\mathcal{L}_{s(i_1)} = (1 - \lambda) \|I'_{s(i_1)} - \tilde{I}_{s(i_1)}\|_1 + \lambda \cdot \text{D-SSIM}(I'_{s(i_1)}, \tilde{I}_{s(i_1)}), \quad (7)$$

where $I'_{s(i_1)}$ denotes the image rendered from the $s(i_1)^{th}$ viewpoint. We set $\lambda = 0.2$ across all experiments and provide comprehensive details of the learning schedule and additional experimental settings in Sec. 4.1.

3.4 Progressive Inpainting

For occlusion-rich, complex scenes, multiple reference views ($r > 1$) are imperative. To solve these challenges, we implement a progressive inpainting approach. Commencing with the initial reference view $s(i_1)$ from the selected views $\mathcal{S} = \{s(i_1), s(i_2), \dots, s(i_r)\}$, we apply Gaussian inpainting as delineated in Sec. 3.3. Subsequent to this, we render the color image, depth map, and associated mask from the next reference view $s(i_2)$. This process is iterated, employing Gaussian inpainting for each successive reference view until the view $s(i_r)$ is addressed. This progressive technique effectively accommodates the complexities, especially for occlusions.

4 Experiments

4.1 Experiments Setup

Training settings To train a diffusion model with broad generalizability for depth inpainting and to facilitate generalized Gaussian inpainting, we train our LDM models using the SceneFlow dataset [62], which comprises FlyingThings and Driving scenes. This dataset offers an extensive collection of over 100,000 frames, each accompanied by ground truth depth, and rendered from a variety of synthetic sequences. During training, masks are randomly generated for each iteration using either a square, random strokes, or a combination of both techniques. We initialize the LDM with pre-trained depth prediction weights sourced from the Marigold [44]. We also tested other pre-trained weight which is presented in the supplementary. The training process spans 200 epochs, with an initial learning rate of 1e-3, which is scheduled to decay after every 50 epochs. Utilizing eight A100 GPUs, the training process is completed within one day.

Evaluation settings We evaluate our method across a variety of datasets, which include forward-facing scenes and the more complex unbounded 360-degree scenes. For the forward-facing datasets, we adhere to the rigorous evaluation settings established by SPIn-NeRF [67]. To further demonstrate the text-guided 3D inpainting capabilities of our method, we also introduce our own captured sequences including large occlusion between objects. The challenging unbounded 3D scenes are taken from the Mip-NeRF [4], featuring large central objects within realistic backgrounds, and the 3DCS, which includes a variety of intricate objects from free-moving camera angles. These datasets are particularly challenging for 3D inpainting. We emphasize that our LDM depth inpainting methods were not trained on any of these datasets. For scene masking, we used masks from SAM-Track [17], dilating them by 9 pixels to reliably remove any undesired parts.

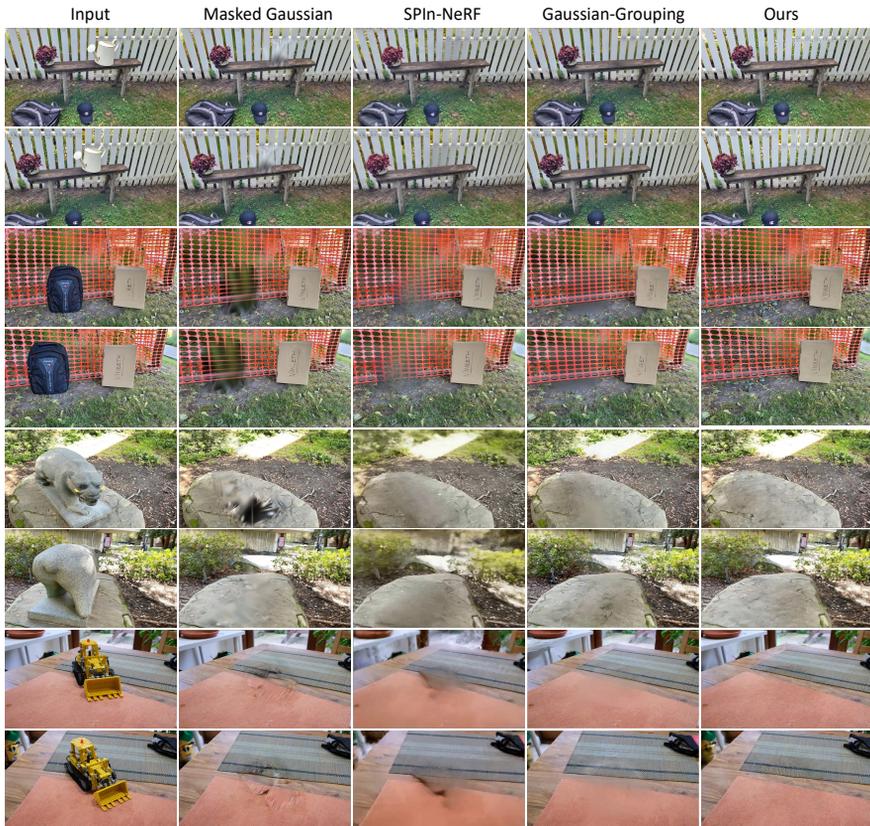


Fig. 3: Qualitative Comparison with Baselines. Zoom in for details. Our method exhibits sharp textures that maintain 3D coherence, whereas baseline approaches often yield details that appear blurred.

Baseline selection Our method has been benchmarked against three distinct baselines, each representing a different approach to inpainting. Spin-NeRF [67] stands out as one of the best NeRF-based methods for 3D inpainting, offering 3D-aware results. For techniques leveraging 3D Gaussians, Gaussian-grouping [106] is the state-of-the-art, building on the InstructNeRF2NeRF [35] framework to incorporate a pre-trained diffusion model for the inpainting task. Lastly, we include a critical baseline where we forgo the depth diffusion inpainting process and instead directly optimize 3D scenes using the inpainted reference image with the aid of Stable Diffusion XL [75].

4.2 Results Comparison

We conduct comparison with baseline methods in several aspects. For our quantitative analysis, we evaluate based on two key metrics: **Image quality** and

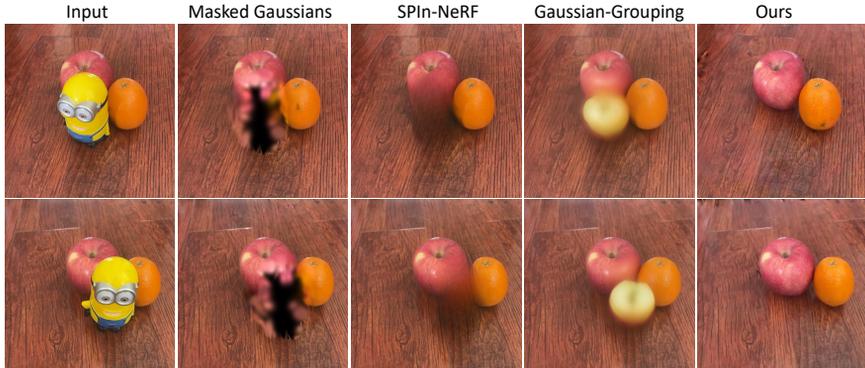


Fig. 4: Qualitative Comparison with Baselines. We delve into more challenging scenarios, including those with multi-object occlusion, where our method uniquely stands out by accurately inpainting the obscured missing segments.

Speed. Image quality is evaluated following the previous works; we report the LPIPS and FID scores for inpainted scenes following the settings in SPIn-NeRF [67]. As detailed in Tab. 1, our method outperforms the baseline techniques in both metrics, achieving the best scores. In terms of speed, our method demonstrates a significant advantage. Thanks to the precision of our initial inpainted points and the efficiency of our fine-tuning process—which necessitates only a minimal number of iterations (around 100)—our approach is considerably faster than baseline methods.

Fig. 3 illustrates a side-by-side comparison of the inpainted results and corresponding novel views generated by our method against those from baseline methods. While baselines are capable of reconstructing the broad outlines of missing regions, they often yield textures that lack sharpness. Our approach, on the other hand, consistently produces fine-detailed textures across all views. Moreover, as shown in Fig. 4, our method can handle more difficult cases which include multi-object occlusion.

	Masked Gaussians	SPIn-NeRF	Gaussians Grouping	Ours
LPIPS ↓	0.594	0.465	0.454	0.421
FID ↓	278.32	156.64	123.48	92.62
Time ↓	20min	5h	20min	40s

Table 1: Quantitative evaluation. We conducted a quantitative evaluation of 3D inpainting techniques on the inpainted areas of held-out views from the SPIn-NeRF dataset. Our method achieves optimal results in perceptual metric (LPIPS) and feature-based statistical distance (FID). Additionally, our method significantly reduces the optimization time compared to previous methods.

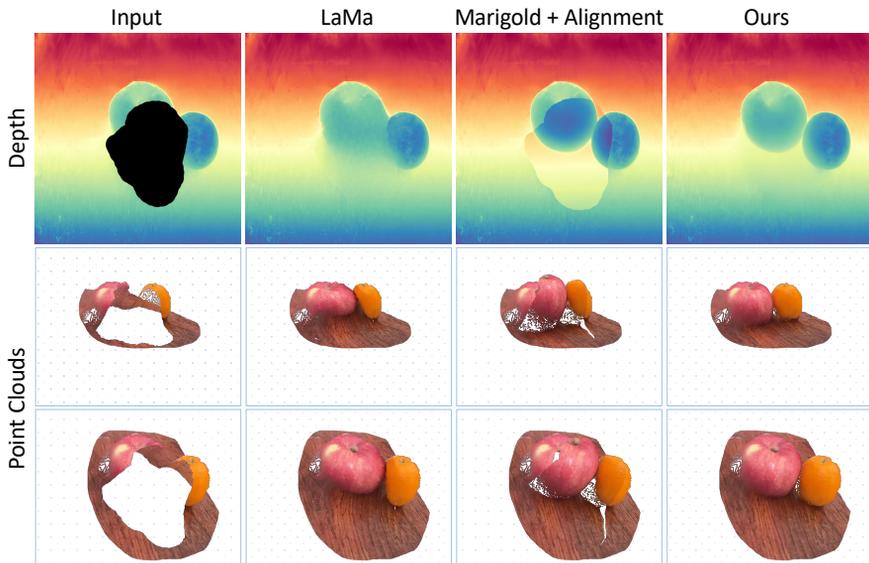


Fig. 5: Ablation study on depth inpainting, we present comparative results against widely-used other baselines, along with the corresponding point cloud visualizations. The comparisons distinctly reveal that our approach successfully inpaints shapes that are correctly aligned with the existing geometry.

Furthermore, we extend our comparative analysis to include additional capabilities of our method relative to the baselines. Owing to our method’s utilization of the reference image, it inherently supports interactive 3D texture inpainting—an operation the baseline methods cannot accommodate. Additionally, our LDM-based approach facilitates object completion in forward-facing scenes. These advanced functionalities are exemplified in the application section of this paper.

4.3 Ablation Study

In order to validate the components within our pipeline, we conducted a series of ablation studies on the key design elements.

Depth Inpainting: A common approach prior to our work involved using image-based inpainting methods, such as LaMa or SDXL Inpainting, for depth inpainting [67,99]. However, because of the domain gap and the model capability, the inpainted results are less accurate as in Fig. 5 and Fig. 1. Another stream involves using monocular depth estimation followed by depth alignment [29]. As depicted in Fig. 5, this method often results in depth discontinuities within the inpainted regions, leading to misalignment with the scene’s original depth. While depth alignment techniques can mitigate this error, significant discrepancies persist.

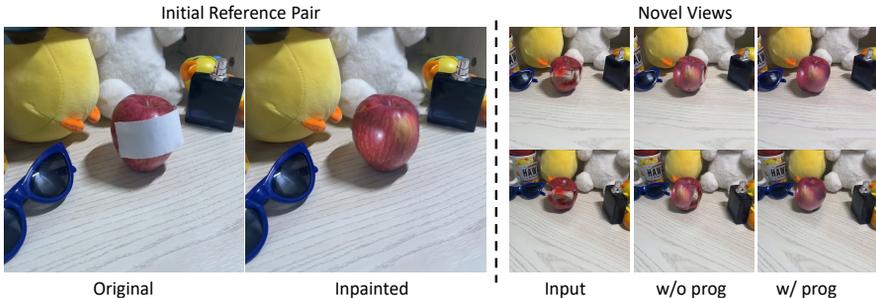


Fig. 6: Ablation study on progressive inpainting. InFusion can adeptly handle inpainting tasks for views that substantially deviate from the initial reference frames.



Fig. 7: User-interactive Texture Inpainting. InFusion allows users to modify the appearance and texture of targeted areas with ease.

Progressive InFusion: Progressive design has a direct impact on performance for difficult cases. As evidenced in our results, augmenting the number of views enhances the handling of occlusions (Fig. 6). Nonetheless, this boost in performance comes at the expense of increased inference time. In simpler scenes, where the task is to remove the outermost object, utilizing a single reference view suffices.

4.4 Applications

To showcase the practical utility of our proposed method, we present two key downstream applications:

Interactive Texture Editing Our framework facilitates user-interactive texture editing within inpainted regions by allowing modifications to the reference image. As illustrated in Fig. 7, users can seamlessly integrate custom text into 3D scenes, enhancing the personalization of the 3D environment.

Object Insertion Leveraging our diffusion-based depth inpainting approach, we enable effortless object insertion within frontal-face scenes, as depicted in Fig. 8. This capability extends to the insertion of user-selected objects into the inpainted 3D scenes, offering a versatile tool for scene customization.

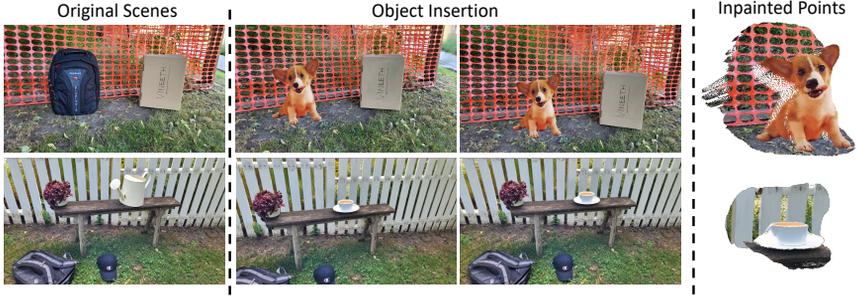


Fig. 8: Object Insertion. Through editing a single image, users are able to project objects into a real three-dimensional scene. This process seamlessly integrates virtual objects into the physical environment, offering an intuitive tool for scene customization.



Fig. 9: Limitations. As the lighting of the surrounding region increasingly differs from the reference, the inpainted area becomes less harmonious with these views. InFusion struggles to adapt inpainted regions to variations in lighting conditions.

4.5 Limitation

While the proposed method achieves impressive results in 3D inpainting, it encounters two main limitations: first, in scenarios with significant lighting changes across various angles, the inpainted sections can struggle to integrate flawlessly with adjacent areas, as highlighted in Fig. 9; second, the method falls short in text-guided inpainting of highly complex objects within 360-degree scenes, limited by the current consistency of inpainting models.

5 Conclusion

In conclusion, our proposed methodology, InFusion, effectively delivers high-quality and efficient inpainting for 3D scenes using Gaussian models. Our evaluations, both quantitative and qualitative, attest to its performance and ease of use. Moreover, we demonstrate that incorporating diffusion priors significantly enhances our depth inpainting model. We are confident that this improved depth inpainting model holds promise for a variety of 3D applications, particularly in

the realm of novel view synthesis. However, our method currently has limitations in handling variations in lighting and reconstructing highly complex structured objects. Despite these challenges, our approach forges a connection between LDM and 3D scene editing. This synergy harbors significant potential for future advancements and optimizations.

References

1. Avrahami, O., Lischinski, D., Fried, O.: Blended diffusion for text-driven editing of natural images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18208–18218 (2022) **2, 5**
2. Ballester, C., Bertalmio, M., Caselles, V., Sapiro, G., Verdera, J.: Filling-in by joint interpolation of vector fields and gray levels. *IEEE Trans. Image Process.* **10**(8), 1200–1211 (2001) **3**
3. Barnes, C., Shechtman, E., Finkelstein, A., Goldman, D.B.: Patchmatch: a randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.* **28**(3), 24 (2009) **3**
4. Barron, J.T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., Srinivasan, P.P.: Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5855–5864 (2021) **1, 9**
5. Bertalmio, M., Sapiro, G., Caselles, V., Ballester, C.: Image inpainting. In: Proceedings of SIGGRAPH. pp. 417–424 (2000) **3, 4**
6. Besic, B., Valada, A.: Dynamic object removal and spatio-temporal RGB-D inpainting via geometry-aware adversarial learning. *IEEE Trans. Intell. Veh.* **7**(2), 170–185 (2022) **4**
7. Betker, J., Goh, G., Jing, L., Brooks, T., Wang, J., Li, L., Ouyang, L., Zhuang, J., Lee, J., Guo, Y., Manassra, W., Dhariwal, P., Chu, C., Jiao, Y., Ramesh, A.: Improving image generation with better captions (2023), <https://cdn.openai.com/papers/dall-e-3.pdf> **4**
8. Bhat, S.F., Birkel, R., Wofk, D., Wonka, P., Müller, M.: Zoedepth: Zero-shot transfer by combining relative and metric depth. *CoRR* **abs/2302.12288** (2023) **4**
9. Blattmann, A., Rombach, R., Ling, H., Dockhorn, T., Kim, S.W., Fidler, S., Kreis, K.: Align your latents: High-resolution video synthesis with latent diffusion models. In: Proceedings of CVPR. pp. 22563–22575 (2023) **2, 4, 5**
10. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of ICCV. pp. 9630–9640 (2021) **4**
11. Cen, J., Fang, J., Yang, C., Xie, L., Zhang, X., Shen, W., Tian, Q.: Segment any 3d gaussians. *arXiv preprint arXiv:2312.00860* (2023) **23**
12. Charatan, D., Li, S., Tagliasacchi, A., Sitzmann, V.: pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. *arXiv preprint arXiv:2312.12337* (2023)
13. Chen, Y., Chen, Z., Zhang, C., Wang, F., Yang, X., Wang, Y., Cai, Z., Yang, L., Liu, H., Lin, G.: Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. *arXiv preprint arXiv:2311.14521* (2023) **2, 4**
14. Chen, Z., Wang, F., Liu, H.: Text-to-3d using gaussian splatting. *arXiv preprint arXiv:2309.16585* (2023) **1**

15. Cheng, H.K., Tai, Y., Tang, C.: Rethinking space-time networks with improved memory coverage for efficient video object segmentation. In: *Advances in NeurIPS 2021*. pp. 11781–11794 (2021) [4](#)
16. Cheng, K.L., Wang, Q., Shi, Z., Zheng, K., Xu, Y., Ouyang, H., Chen, Q., Shen, Y.: Learning naturally aggregated appearance for efficient 3d editing. *CoRR* **abs/2312.06657** (2023). <https://doi.org/10.48550/ARXIV.2312.06657>, <https://doi.org/10.48550/arXiv.2312.06657> [4](#)
17. Cheng, Y., Li, L., Xu, Y., Li, X., Yang, Z., Wang, W., Yang, Y.: Segment and track anything. *arXiv preprint arXiv:2305.06558* (2023) [9](#)
18. Chu, T., Chen, J., Sun, J., Lian, S., Wang, Z., Zuo, Z., Zhao, L., Xing, W., Lu, D.: Rethinking fast fourier convolution in image inpainting. In: *Proceedings of ICCV*. pp. 23138–23148 (2023) [3](#)
19. Chung, J., Lee, S., Nam, H., Lee, J., Lee, K.M.: Luciddreamer: Domain-free generation of 3d gaussian splatting scenes. *arXiv preprint arXiv:2311.13384* (2023) [5](#)
20. Dai, A., Diller, C., Nießner, M.: SG-NN: sparse generative neural networks for self-supervised scene completion of RGB-D scans. In: *Proceedings of CVPR*. pp. 846–855 (2020) [4](#)
21. Dai, A., Qi, C.R., Nießner, M.: Shape completion using 3d-encoder-predictor cnns and shape synthesis. In: *Proceedings of CVPR*. pp. 6545–6554 (2017) [4](#)
22. Dai, A., Ritchie, D., Bokeloh, M., Reed, S., Sturm, J., Nießner, M.: Scancomplete: Large-scale scene completion and semantic segmentation for 3d scans. In: *Proceedings of CVPR*. pp. 4578–4587 (2018) [4](#)
23. Darabi, S., Shechtman, E., Barnes, C., Goldman, D.B., Sen, P.: Image melding: combining inconsistent images using patch-based synthesis. *ACM Trans. Graph.* **31**(4), 82:1–82:10 (2012) [3](#)
24. Dong, Q., Cao, C., Fu, Y.: Incremental transformer structure enhanced image inpainting with masking positional encoding. In: *Proceedings of CVPR*. pp. 11348–11358 (2022) [3](#)
25. Dou, B., Zhang, T., Ma, Y., Wang, Z., Yuan, Z.: Cosseggaussians: Compact and swift scene segmenting 3d gaussians. *arXiv preprint arXiv:2401.05925* (2024) [23](#)
26. Duan, Y., Zhu, Z., Guo, X.: Diffusiondepth: Diffusion denoising approach for monocular depth estimation. *CoRR* **abs/2303.05021** (2023) [4](#)
27. Efros, A.A., Leung, T.K.: Texture synthesis by non-parametric sampling. In: *Proceedings of ICCV*. pp. 1033–1038 (1999) [3](#)
28. Eigen, D., Fergus, R.: Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In: *Proceedings of ICCV*. pp. 2650–2658 (2015) [4](#)
29. Fang, J., Wang, J., Zhang, X., Xie, L., Tian, Q.: Gaussianeditor: Editing 3d gaussians delicately with text instructions. *arXiv preprint arXiv:2311.16037* (2023) [2](#), [12](#)
30. Fujii, R., Hachiuma, R., Saito, H.: RGB-D image inpainting using generative adversarial network with a late fusion approach. In: *Proceedings of AVR*. vol. 12242, pp. 440–451 (2020) [4](#)
31. Gao, C., Saraf, A., Huang, J., Kopf, J.: Flow-edge guided video completion. In: *Proceedings of ECCV*. vol. 12357, pp. 713–729 (2020) [4](#)
32. Godard, C., Aodha, O.M., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: *Proceedings of CVPR*. pp. 6602–6611 (2017) [4](#)

33. Granados, M., Kim, K.I., Tompkin, J., Kautz, J., Theobalt, C.: Background inpainting for videos with dynamic objects and a free-moving camera. In: Proceedings of ECCV. vol. 7572, pp. 682–695 (2012) [3](#)
34. Han, X., Li, Z., Huang, H., Kalogerakis, E., Yu, Y.: High-resolution shape completion using deep neural networks for global structure and local geometry inference. In: Proceedings of ICCV. pp. 85–93 (2017) [4](#)
35. Haque, A., Tancik, M., Efros, A.A., Holynski, A., Kanazawa, A.: Instruct-nerf2nerf: Editing 3d scenes with instructions. arXiv preprint arXiv:2303.12789 (2023) [4](#), [10](#)
36. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. In: Advances in NeurIPS (2020) [4](#), [7](#)
37. Hu, X., Wang, Y., Fan, L., Fan, J., Peng, J., Lei, Z., Li, Q., Zhang, Z.: Semantic anything in 3d gaussians. arXiv preprint arXiv:2401.17857 (2024) [23](#)
38. Huang, J., Kang, S.B., Ahuja, N., Kopf, J.: Temporally coherent completion of dynamic video. ACM Trans. Graph. **35**(6), 196:1–196:11 (2016) [3](#)
39. Huang, J., Yu, H.: Point’n move: Interactive scene object manipulation on gaussian splatting radiance fields. CoRR **abs/2311.16737** (2023) [4](#)
40. Ji, Y., Chen, Z., Xie, E., Hong, L., Liu, X., Liu, Z., Lu, T., Li, Z., Luo, P.: DDP: diffusion model for dense visual prediction. In: Proceedings of ICCV. pp. 21684–21695 (2023) [4](#)
41. Jiang, Y., Yu, C., Xie, T., Li, X., Feng, Y., Wang, H., Li, M., Lau, H.Y.K., Gao, F., Yang, Y., Jiang, C.: VR-GS: A physical dynamics-aware interactive gaussian splatting system in virtual reality. CoRR **abs/2401.16663** (2024) [4](#)
42. Kazhdan, M.M., Bolitho, M., Hoppe, H.: Poisson surface reconstruction. In: Proceedings of SGP. pp. 61–70 (2006) [4](#)
43. Kazhdan, M.M., Hoppe, H.: Screened poisson surface reconstruction. ACM Trans. Graph. **32**(3), 29:1–29:13 (2013) [4](#)
44. Ke, B., Obukhov, A., Huang, S., Metzger, N., Dautt, R.C., Schindler, K.: Repurposing diffusion-based image generators for monocular depth estimation. CoRR **abs/2312.02145** (2023) [2](#), [3](#), [4](#), [5](#), [7](#), [9](#), [22](#)
45. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. ACM Transactions on Graphics **42**(4) (2023) [1](#)
46. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. ACM Trans. Graph. **42**(4), 139:1–139:14 (2023) [4](#), [5](#)
47. Kerr, J., Kim, C.M., Goldberg, K., Kanazawa, A., Tancik, M.: LERF: language embedded radiance fields. In: Proceedings of ICCV. pp. 19672–19682 (2023) [4](#)
48. Ko, K., Kim, C.: Continuously masked transformer for image inpainting. In: Proceedings of ICCV. pp. 13123–13132 (2023) [3](#)
49. Kobayashi, S., Matsumoto, E., Sitzmann, V.: Decomposing nerf for editing via feature field distillation. In: Advances in NeurIPS (2022) [4](#)
50. Lan, K., Li, H., Shi, H., Wu, W., Liao, Y., Wang, L., Zhou, P.: 2d-guided 3d gaussian segmentation. arXiv preprint arXiv:2312.16047 (2023) [23](#)
51. Lee, E., Yoo, J., Yang, Y., Baik, S., Kim, T.H.: Semantic-aware dynamic parameter for video inpainting transformer. In: Proceedings of ICCV. pp. 12903–12912 (2023) [4](#)
52. Li, W., Lin, Z., Zhou, K., Qi, L., Wang, Y., Jia, J.: MAT: mask-aware transformer for large hole image inpainting. In: Proceedings of CVPR. pp. 10748–10758 (2022) [3](#)

53. Li, Z., Lu, C., Qin, J., Guo, C., Cheng, M.: Towards an end-to-end framework for flow-guided video inpainting. In: Proceedings of CVPR. pp. 17541–17550 (2022) [4](#)
54. Li, Z., Snavely, N.: Megadepth: Learning single-view depth prediction from internet photos. In: Proceedings of CVPR. pp. 2041–2050 (2018) [4](#)
55. Liu, H., Shen, I., Chen, B.: Nerf-in: Free-form nerf inpainting with RGB-D priors. CoRR [abs/2206.04901](#) (2022) [4](#)
56. Liu, R., Deng, H., Huang, Y., Shi, X., Lu, L., Sun, W., Wang, X., Dai, J., Li, H.: Fuseformer: Fusing fine-grained information in transformers for video inpainting. In: Proceedings of ICCV. pp. 14020–14029 (2021) [4](#)
57. Liu, W., Chen, X., Yang, J., Wu, Q.: Robust color guided depth map restoration. IEEE Trans. Image Process. **26**(1), 315–327 (2017) [4](#)
58. Liu, Y., Hu, B., Huang, J., Tai, Y., Tang, C.: Instance neural radiance field. In: Proceedings of ICCV. pp. 787–796 (2023) [4](#)
59. Lu, S., Ren, X., Liu, F.: Depth enhancement via low-rank matrix completion. In: Proceedings of CVPR. pp. 3390–3397 (2014) [4](#)
60. Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., Gool, L.V.: Repaint: Inpainting using denoising diffusion probabilistic models. In: Proceedings of CVPR. pp. 11451–11461 (2022) [3](#)
61. Makarov, I., Borisenko, G.: Depth inpainting via vision transformer. In: Proceedings of ISMAR. pp. 286–291 (2021) [4](#)
62. Mayer, N., Ilg, E., Hausser, P., Fischer, P., Cremers, D., Dosovitskiy, A., Brox, T.: A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4040–4048 (2016) [2, 9](#)
63. Menze, M., Heipke, C., Geiger, A.: Object scene flow. ISPRS Journal of Photogrammetry and Remote Sensing **140**, 60–76 (2018) [2](#)
64. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: Proceedings of ECCV. vol. 12346, pp. 405–421 (2020) [4](#)
65. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM **65**(1), 99–106 (2021) [1](#)
66. Mirzaei, A., Aumentado-Armstrong, T., Brubaker, M.A., Kelly, J., Levinshstein, A., Derpanis, K.G., Gilitschenski, I.: Reference-guided controllable inpainting of neural radiance fields. In: Proceedings of ICCV. pp. 17769–17779 (2023) [4](#)
67. Mirzaei, A., Aumentado-Armstrong, T., Derpanis, K.G., Kelly, J., Brubaker, M.A., Gilitschenski, I., Levinshstein, A.: Spin-nerf: Multiview segmentation and perceptual inpainting with neural radiance fields. In: Proceedings of CVPR. pp. 20669–20679 (2023) [3, 4, 9, 10, 11, 12](#)
68. Mirzaei, A., Kant, Y., Kelly, J., Gilitschenski, I.: Laterf: Label and text driven object radiance fields. In: Proceedings of ECCV. pp. 20–36 (2022) [4](#)
69. Newson, A., Almansa, A., Fradet, M., Gousseau, Y., Pérez, P.: Towards fast, generic video inpainting. In: Proceedings of CVMP. pp. 7:1–7:8 (2013) [3](#)
70. Newson, A., Almansa, A., Fradet, M., Gousseau, Y., Pérez, P.: Video inpainting of complex scenes. Siam journal on imaging sciences **7**(4), 1993–2019 (2014) [3](#)
71. Ouyang, H., Heal, K., Lombardi, S., Sun, T.: Text2immersion: Generative immersive scene with 3d gaussians. arXiv preprint arXiv:2312.09242 (2023) [5](#)
72. Ouyang, H., Wang, Q., Xiao, Y., Bai, Q., Zhang, J., Zheng, K., Zhou, X., Chen, Q., Shen, Y.: Codef: Content deformation fields for temporally consistent video processing. CoRR [abs/2308.07926](#) (2023) [3](#)

73. Pang, H., Hua, B., Yeung, S.: Locally stylized neural radiance fields. CoRR **abs/2309.10684** (2023) [4](#)
74. Park, J.J., Florence, P.R., Straub, J., Newcombe, R.A., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. In: Proceedings of CVPR. pp. 165–174 (2019) [4](#)
75. Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., Rombach, R.: SDXL: improving latent diffusion models for high-resolution image synthesis. CoRR **abs/2307.01952** (2023) [2](#), [3](#), [4](#), [5](#), [7](#), [10](#), [24](#)
76. Quan, W., Chen, J., Liu, Y., Yan, D., Wonka, P.: Deep learning-based image and video inpainting: A survey. CoRR **abs/2401.03395** (2024) [3](#)
77. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: Proceedings of ICML. pp. 8748–8763 (2021) [4](#)
78. Ranftl, R., Lasinger, K., Hafner, D., Schindler, K., Koltun, V.: Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. IEEE Trans. Pattern Anal. Mach. Intell. **44**(3), 1623–1637 (2022) [4](#)
79. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022) [2](#), [3](#), [4](#), [5](#), [6](#), [22](#)
80. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, S.K.S., Lopes, R.G., Ayan, B.K., Salimans, T., Ho, J., Fleet, D.J., Norouzi, M.: Photorealistic text-to-image diffusion models with deep language understanding. In: Advances in NeurIPS (2022) [4](#)
81. Sargsyan, A., Navasardyan, S., Xu, X., Shi, H.: MI-GAN: A simple baseline for image inpainting on mobile devices. In: Proceedings of ICCV. pp. 7301–7311 (2023) [3](#)
82. Saxena, S., Herrmann, C., Hur, J., Kar, A., Norouzi, M., Sun, D., Fleet, D.J.: The surprising effectiveness of diffusion models for optical flow and monocular depth estimation. In: Advances in NeurIPS (2023) [4](#)
83. Saxena, S., Kar, A., Norouzi, M., Fleet, D.J.: Monocular depth estimation using diffusion models. CoRR **abs/2302.14816** (2023) [4](#)
84. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5b: An open large-scale dataset for training next generation image-text models. Advances in Neural Information Processing Systems **35**, 25278–25294 (2022) [2](#)
85. Shih, M., Su, S., Kopf, J., Huang, J.: 3d photography using context-aware layered depth inpainting. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020. pp. 8025–8035 (2020) [4](#)
86. Shih, T.K., Tang, N.C., Hwang, J.: Exemplar-based video inpainting without ghost shadow artifacts by maintaining temporal continuity. IEEE Trans. Circuits Syst. Video Technol. **19**(3), 347–360 (2009) [3](#)
87. Siddiqui, Y., Porzi, L., Bulò, S.R., Müller, N., Nießner, M., Dai, A., Kotschieder, P.: Panoptic lifting for 3d scene understanding with neural fields. In: Proceedings of CVPR. pp. 9043–9052 (2023) [4](#)
88. Sohl-Dickstein, J., Weiss, E.A., Maheswaranathan, N., Ganguli, S.: Deep unsupervised learning using nonequilibrium thermodynamics. In: Proceedings of ICML. pp. 2256–2265 (2015) [4](#)

89. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. In: Proceedings of ICLR (2021) [4](#), [8](#)
90. Song, S., Yu, F., Zeng, A., Chang, A.X., Savva, M., Funkhouser, T.A.: Semantic scene completion from a single depth image. In: Proceedings of CVPR. pp. 190–198 (2017) [4](#)
91. Strobel, M., Diebold, J., Cremers, D.: Flow and color inpainting for video completion. In: Pattern Recognition - 36th German Conference, GCPR 2014, Münster, Germany, September 2-5, 2014, Proceedings. vol. 8753, pp. 293–304 (2014) [3](#)
92. Suvorov, R., Logacheva, E., Mashikhin, A., Remizova, A., Ashukha, A., Silvestrov, A., Kong, N., Goka, H., Park, K., Lempitsky, V.: Resolution-robust large mask inpainting with fourier convolutions. In: Proceedings of WACV. pp. 3172–3182. IEEE (2022) [2](#), [3](#), [4](#)
93. Tang, J., Ren, J., Zhou, H., Liu, Z., Zeng, G.: Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. arXiv preprint arXiv:2309.16653 (2023) [1](#)
94. Tschumperlé, D., Deriche, R.: Vector-valued image regularization with pdes: A common framework for different applications. IEEE Trans. Pattern Anal. Mach. Intell. **27**(4), 506–517 (2005) [3](#)
95. Wang, D., Zhang, T., Abboud, A., Süssstrunk, S.: Inpaintnerf360: Text-guided 3d inpainting on unbounded neural radiance fields. CoRR [abs/2305.15094](#) (2023) [4](#)
96. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. arXiv preprint arXiv:2106.10689 (2021) [1](#)
97. Wang, W., Huang, Q., You, S., Yang, C., Neumann, U.: Shape inpainting using 3d generative adversarial network and recurrent convolutional networks. In: Proceedings of ICCV. pp. 2317–2325 (2017) [4](#)
98. Weber, E., Holynski, A., Jampani, V., Saxena, S., Snavely, N., Kar, A., Kanazawa, A.: Nerfiller: Completing scenes via generative 3d inpainting. CoRR [abs/2312.04560](#) (2023) [4](#)
99. Weder, S., Garcia-Hernando, G., Monszpart, Á., Pollefeys, M., Brostow, G.J., Firman, M., Vicente, S.: Removing objects from neural radiance fields. In: Proceedings of CVPR. pp. 16528–16538 (2023) [4](#), [12](#)
100. Wei, F., Funkhouser, T.A., Rusinkiewicz, S.: Clutter detection and removal in 3d scenes with view-consistent inpainting. In: Proceedings of ICCV. pp. 18085–18095 (2023) [4](#)
101. Wexler, Y., Shechtman, E., Irani, M.: Space-time completion of video. IEEE Trans. Pattern Anal. Mach. Intell. **29**(3), 463–476 (2007) [3](#)
102. Wu, G., Yi, T., Fang, J., Xie, L., Zhang, X., Wei, W., Liu, W., Tian, Q., Wang, X.: 4d gaussian splatting for real-time dynamic scene rendering. arXiv preprint arXiv:2310.08528 (2023) [1](#)
103. Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: Proceedings of CVPR. pp. 1912–1920 (2015) [4](#)
104. Xu, R., Li, X., Zhou, B., Loy, C.C.: Deep flow-guided video inpainting. In: Proceedings of CVPR. pp. 3723–3732 (2019) [3](#), [4](#)
105. Yang, L., Kang, B., Huang, Z., Xu, X., Feng, J., Zhao, H.: Depth anything: Unleashing the power of large-scale unlabeled data. CoRR [abs/2401.10891](#) (2024) [4](#), [24](#)

106. Ye, M., Danelljan, M., Yu, F., Ke, L.: Gaussian grouping: Segment and edit anything in 3d scenes. arXiv preprint arXiv:2312.00732 (2023) [2](#), [3](#), [4](#), [10](#), [23](#)
107. Yi, T., Fang, J., Wu, G., Xie, L., Zhang, X., Liu, W., Tian, Q., Wang, X.: Gaussiandreamer: Fast generation from text to 3d gaussian splatting with point cloud priors. arXiv preprint arXiv:2310.08529 (2023) [1](#)
108. Yu, T., Feng, R., Feng, R., Liu, J., Jin, X., Zeng, W., Chen, Z.: Inpaint anything: Segment anything meets image inpainting. CoRR **abs/2304.06790** (2023) [3](#)
109. Zeng, Y., Fu, J., Chao, H.: Learning joint spatial-temporal transformations for video inpainting. In: Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVI 16. pp. 528–543. Springer (2020) [3](#)
110. Zhang, K., Kolkin, N.I., Bi, S., Luan, F., Xu, Z., Shechtman, E., Snavely, N.: ARF: artistic radiance fields. In: Proceedings of ECCV (2022) [4](#)
111. Zhang, K., Fu, J., Liu, D.: Flow-guided transformer for video inpainting. In: Proceedings of ECCV. vol. 13678, pp. 74–90 (2022) [4](#)
112. Zhang, K., Fu, J., Liu, D.: Inertia-guided flow completion and style fusion for video inpainting. In: Proceedings of CVPR. pp. 5972–5981 (2022) [4](#)
113. Zhang, Y., Funkhouser, T.A.: Deep depth completion of a single RGB-D image. In: Proceedings of CVPR. pp. 175–185 (2018) [4](#)
114. Zhang, Y., He, Z., Xing, J., Yao, X., Jia, J.: Ref-npr: Reference-based non-photorealistic radiance fields for controllable scene stylization. In: Proceedings of CVPR (2023) [4](#)
115. Zhang, Y., Scargill, T., Vaishnav, A., Premsankar, G., Francesco, M.D., Gorlatova, M.: Indepth: Real-time depth inpainting for mobile augmented reality. Proc. ACM Interact. Mob. Wearable Ubiquitous Technol. **6**(1), 37:1–37:25 (2022) [4](#)
116. Zhao, W., Rao, Y., Liu, Z., Liu, B., Zhou, J., Lu, J.: Unleashing text-to-image diffusion models for visual perception. In: Proceedings of ICCV. pp. 5706–5716 (2023) [4](#)
117. Zheng, W., Xu, C., Xu, X., Liu, W., He, S.: CIRI: curricular inactivation for residue-aware one-shot video inpainting. In: Proceedings of ICCV. pp. 12966–12976 (2023) [4](#)
118. Zhou, S., Li, C., Chan, K.C.K., Loy, C.C.: Propainter: Improving propagation and transformer for video inpainting. In: Proceedings of ICCV. pp. 10443–10452 (2023) [4](#)
119. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: Proceedings of CVPR. pp. 6612–6619 (2017) [4](#)
120. Zhuang, J., Wang, C., Lin, L., Liu, L., Li, G.: Dreameditor: Text-driven 3d scene editing with neural fields. In: Kim, J., Lin, M.C., Bickel, B. (eds.) Proceedings of SIGGRAPH Asia. pp. 26:1–26:10 (2023) [4](#)

The appendix are structured as follows: We begin with a detailed description of our implementation, including the specifics of our training configurations and the outlier removal process for the point cloud. Subsequently, we undertake an in-depth examination of various issues discussed within the paper. To conclude, we provide a broader set of results, encompassing a wider array of scenes and viewpoints.

1 Implementation Details

1.1 Training details

The Depth inpainting model is initialized with the Marigold [44] weights. The architecture of the neural network is consistent with that of Stable Diffusion v1.5 [79], with the exception of the first convolutional layer. Moreover, during both training and inference phases, the input to the text encoder is persistently an empty string. The UNet has 9 additional input channels (4 for the encoded masked-depth, 4 for the guided encoded image and 1 for the mask itself) whose weights were zero-initialized. During training, we generate synthetic masks and, in 30% mask everything. In the context of data processing, we maintain the original aspect ratio of the images during both the training and inference stages, resizing them to a maximum resolution of 768 pixels on the longest side.

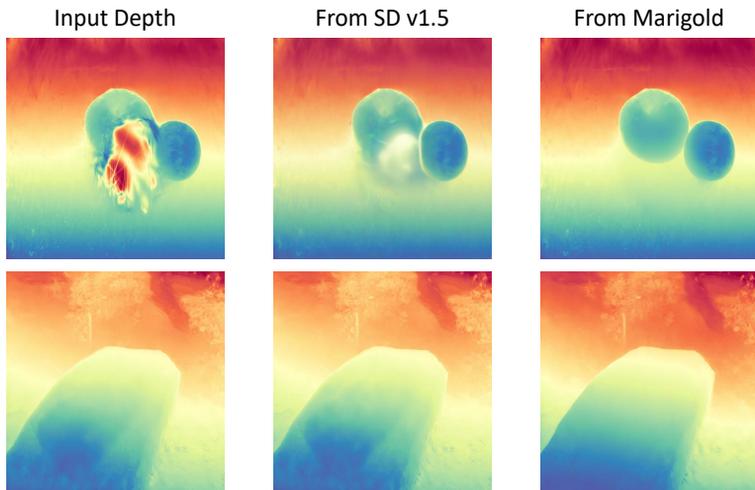


Fig. 1: Analysis on Pre-trained Weights.

1.2 Outliers removal

we unproject depth map and reference image from image space to 3D coordinates to form a colored point cloud. Before this point cloud is merged into original 3D Gaussian point cloud, we need process outliers to improve rendered

image quality. To eliminate Gaussian outliers along the edges of the mask, we initially construct a KDTree from the unprojected point cloud. Subsequently, this KDTree is employed to locate the nearest points within the original point cloud, returning points from the original cloud that are within a specified distance threshold. Subsequently, we utilize the *'remove_radius_outlier'* method from the point cloud data (pcd) library to identify points in the original point cloud that have an insufficient number of neighbors within a specified radius. An intersection of these points and the similar points previously determined using a KDTree is performed, thereby efficiently removing Gaussian outliers at the edges of the mask. Additionally, there are various Gaussian segmentation [11, 25, 37, 50, 106] techniques that can be employed for outlier removal, taking advantage of the explicit properties of Gaussian models. Nevertheless, these are not the focal point of the present study and will not be deliberated here.

2 Analysis

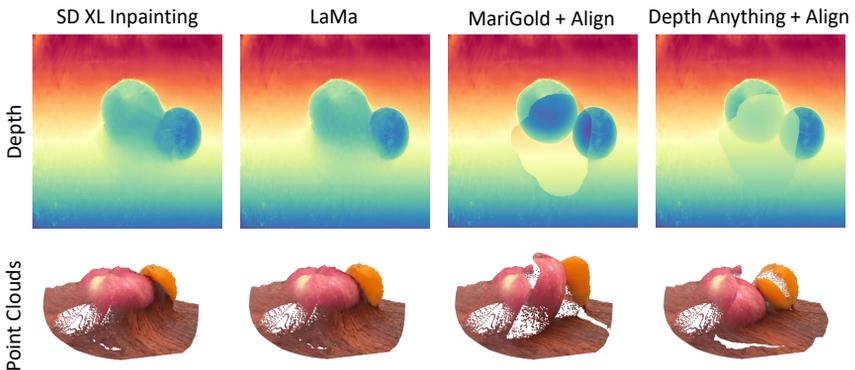


Fig. 2: Analysis on Depth Inpainting. It is evident that the image-based inpainting models, lacking proper guidance, fail to adequately complete the geometric details. Regarding the monocular estimation methods, while a depth alignment method is implemented, they often lead to discontinuities within the inpainted regions

2.1 Analysis on pre-trained weights

For the task of depth completion, we employed two distinct sets of initial weights: one derived from Marigold and the other based on Stable Diffusion v1.5. As demonstrated in Fig. 1, we display the results under conditions of equivalent data volume and identical training epochs. It is discerned that models initialized with weights from Stable Diffusion v1.5 encountered greater challenges in mastering the depth completion task, a difficulty that was particularly pronounced in complex scenes. In contrast, models that began with Marigold weights exhibited superior proficiency in completing depth, due to prior training on depth maps

that reduced the gap between the RGB and depth domains. Following the same training regimen, these models demonstrated an enhanced ability for depth completion and achieved better alignment with the input images.

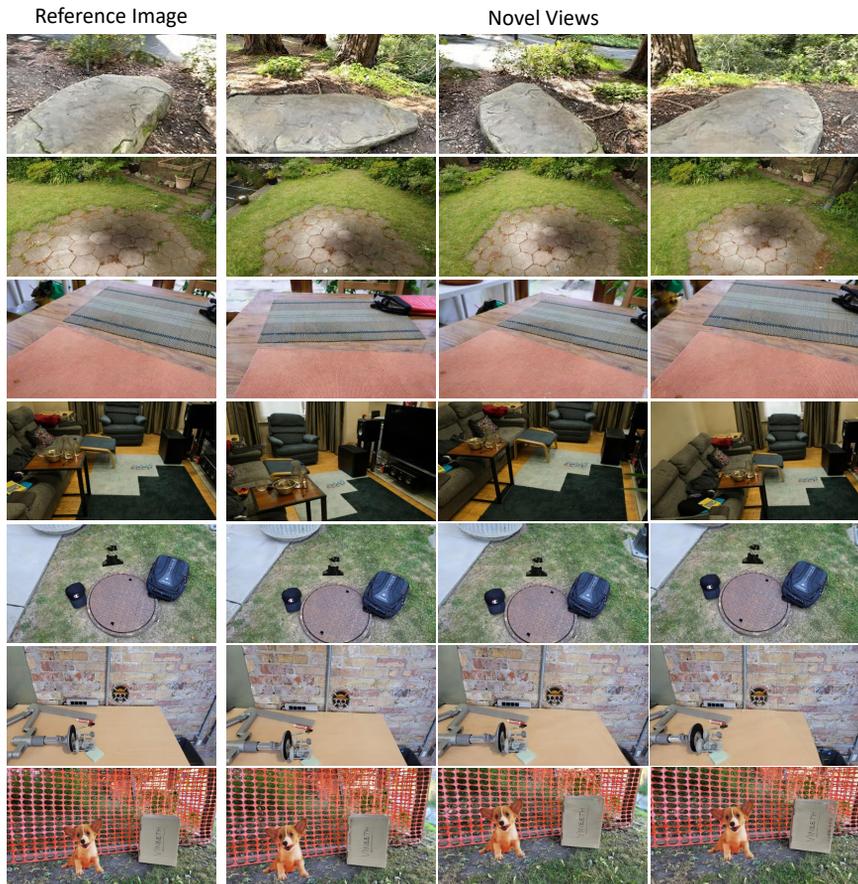


Fig. 3: Qualitative Results. Zoom in for details. Our method exhibits sharp textures that maintain 3D coherence. We respectfully invite you to view the video featured on the webpage within our supplementary materials

2.2 Analysis on depth inpainting

We include feature additional results, comparing our method with various cutting-edge baselines, such as SD XL inpainting [75] and DepthAnything [105], with a focus on alignment accuracy. As shown in Fig. 2, while SD XL inpainting

yields visually appealing results in the RGB domain, a closer inspection of the reprojected point clouds reveals noticeable inaccuracies, akin to those observed in LaMa. Similarly, DepthAnything struggles with discontinuities, leading to a pronounced gap between inpainted areas and their adjacent regions, much like the issues seen with MariGold. Consequently, our learned depth inpainting is critical in securing high-fidelity results.

3 More Results

As shown in Fig. 3, we present the single reference images for several scenes , along with multiple novel views, to validate the robust 3D consistency achieved by InFusion. Additionally, we have consolidated all scenes into a webpage included in our supplementary materials and extend an invitation for you to view them.