

3D Gaussian Inpainting with Depth-Guided Cross-View Consistency

Sheng-Yu Huang^{1,†}, Zi-Ting Chou², Yu-Chiang Frank Wang^{1,2,‡},

¹ Graduate Institute of Communication Engineering, National Taiwan University ² NVIDIA, Taiwan

[†]f08942095@ntu.edu.tw [‡]frankwang@nvidia.com

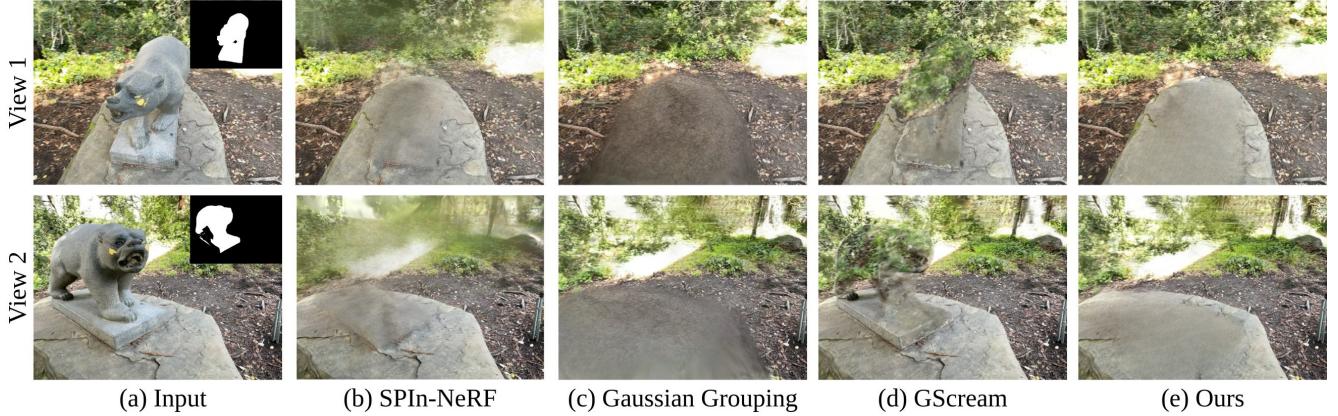


Figure 1. Given multi-view images of a scene and the object masks describing the object to be removed, our 3D Gaussian Inpainting with Depth-Guided Cross-View Consistency produces high-fidelity cross-view inpainting results. Compared with current state-of-the-arts such as SPIIn-NeRF [25], Gaussian Grouping [44], and GScream [38], inpainting results of our method not only preserve visible background contents but also exhibit satisfactory consistency across camera views.

Abstract

When performing 3D inpainting using novel-view rendering methods like Neural Radiance Field (NeRF) or 3D Gaussian Splatting (3DGS), how to achieve texture and geometry consistency across camera views has been a challenge. In this paper, we propose a framework of 3D Gaussian Inpainting with Depth-Guided Cross-View Consistency (3DGIC) for cross-view consistent 3D inpainting. Guided by the rendered depth information from each training view, our 3DGIC exploits background pixels visible across different views for updating the inpainting mask, allowing us to refine the 3DGS for inpainting purposes. Through extensive experiments on benchmark datasets, we confirm that our 3DGIC outperforms current state-of-the-art 3D inpainting methods quantitatively and qualitatively.

1. Introduction

Novel view synthesis for 3D scenes plays a vital role in 3D reconstruction and scene understanding. Recent advancements, such as Neural Radiance Fields (NeRF) [1, 5, 10, 22, 23, 27, 31, 35, 46] and 3D Gaussian Splatting (3DGS) [6, 15, 30, 44, 47], enable high-fidelity novel views

by modeling volumetric properties. However, practical VR/AR applications [3, 21] require more than reconstruction: they need editing capabilities that these methods do not fully address. Among editing challenges, object removal and inpainting [18, 38] are particularly difficult, as direct removal creates visible holes, compromising visual quality. While 2D inpainting [8, 20, 29, 33, 36, 41–43] across multiple views is possible, maintaining consistency remains problematic, leading to artifacts and reduced fidelity. Thus, achieving seamless, multi-view-consistent inpainting for 3D scenes is still an open challenge.

As a pioneering work in 3D scene inpainting, SPIIn-NeRF [25] proposes to use a pre-trained segmentation network [12] to generate plausible 2D inpaint masks for multi-view images with sparse human annotations of the object to be removed. However, as noted in subsequent research [18, 38], SPIIn-NeRF and similar approaches [40, 45] rely heavily on 2D inpainting of multiple views separately, which hinders the cross-view consistency of the 3D inpainting results. To ensure cross-view consistency, RefNeRF [24] projects the inpainted image from a specific reference view onto other views using depth-guided projection, thereby ensuring more consistent inpainting results across views. Despite these advancements, these methods still require human-

annotated 2D masks or sparse annotations to delineate the objects to be removed and the regions to be inpainted, making the process labor-intensive and limiting the scalability and practicality of these techniques.

To reduce the need for human annotation for obtaining inpainting masks, recent methods [44, 45] tend to leverage the Segment Anything Model (SAM) [17] models with NeRF or 3DGS to obtain 2D inpainting masks for multi-view images directly. Although these methods ease the requirement of human annotations for inpainting masks, they still rely on 2D inpainting results for different views as supervision, limiting the multi-view consistency of the inpainted 3D representations. To alleviate this limitation, some approaches [7, 18, 19, 26, 38, 39] attempt to build a cross-view consistent 3D inpainting method on top of the 2D inpainting mask obtained from SAM. By either leveraging 2D diffusion models as perceptual guidance for the inpainted region [7, 18, 39] or ensuring feature consistency of corresponding pixels across different views [38], these methods are able to produce more consistent 3D inpainting results without the requirement of human-annotated 2D inpainting masks. Nevertheless, most of the aforementioned methods rely on the provided per-scene 2D inpainting masks (either from human annotation or from SAM) for each view, which can include areas visible in other views, as mentioned in [44]. As a result, the inpainted content within this area might be inconsistent across camera views, producing artifacts in the reconstructed 3D scene.

In this paper, we propose a 3D Gaussian Inpainting with Depth-Guided Cross-View Consistency (3DGIC) to optimize the 3DGS model while achieving multi-view consistent and high-fidelity 3D inpainting with depth-guided inpainting masks to locate the inpainting region. Given a set of images of a scene with corresponding camera views and the object masks indicating an unwanted object in the scene (obtained from SAM [17], for example), our 3DGIC conducts the process of *Inferring Depth-Guided Inpainting Masks* to consider depth information from all training views and refine the inpainting mask by discovering background pixels from different views. The refined inpainting masks are then used to provide a joint update of inpainting results and the underlying 3DGS model via *3D Inpainting with Cross-View Consistency*. Through experiments on real-world datasets, we quantitatively and qualitatively demonstrate that our 3DGIC performs favorably against state-of-the-art NeRF/3DGS-based inpainting methods by achieving better fidelity and multi-view consistency.

The key contributions of our approach are as follows:

- We propose a 3D Gaussian Inpainting with Depth-Guided Cross-View Consistency (3DGIC), achieving multi-view consistent 3D inpainting results with high fidelity.
- By inferring Depth-Guided Inpainting Masks, the region to be inpainted is properly obtained by considering depth

information across different views, allowing us to guide the inpainting process for 3DGS.

- Based on the 2D inpaintings from a chosen reference view, our Inpainting-guided 3DGS Refinement optimizes new Gaussians of the object-removed scene by ensuring cross-view consistent inpainting results.

2. Related Works

2.1. 3D Representations for Novel View Synthesis

Novel view synthesis is a widely studied topic in 3D computer vision. Neural Radiance Field (NeRF) [23], a pioneer in this field, effectively models scenes using multi-view images. However, as noted in [9], the original NeRF requires extensive training time—from hours to days—and relies on numerous images. To address these issues, many subsequent works [10, 27, 35, 46] have emerged. Methods like Instant NGP [27] and DVGO [35] reduce training time to minutes by balancing speed and memory through hash encoding and voxel encoding. Recently, the introduction of 3D Gaussian Splatting (3DGS) [15] brings a fundamental revolution to this area. Different from NeRF and its variants, which model a 3D scene as an implicit representation, 3DGS models a 3D scene as a composition of numerous 3D Gaussians, with each Gaussian parameterized by its three-dimensional centroid, standard deviations, orientations, opacity, and color features. By modeling a 3D scene as such an explicit representation, one is able to render the 2D images of the modeled scene via rasterization with an incredible 100 fps, whereas the fastest NeRF-based approach ([10, 27]) only achieves around 10 fps. As a result, we chose 3DGS as our backbone representation over NeRF in this paper due to its fast rendering property, making our approach more applicable in the real world.

2.2. 3D Scene Inpainting

In the context of 3D scene inpainting, SPIn-NeRF [25] emerges as one of the earliest approaches addressing the challenges of multi-view consistency. It uses pre-trained segmentation networks to generate plausible inpainting masks for multi-view images, requiring sparse user annotations to indicate the unwanted object. These annotations are propagated across views, and a modified Neural Radiance Field (NeRF) model is used to inpaint the masked regions. Although effective, this approach is heavily dependent on human intervention and lacks the ability to automate the mask generation process, thus limiting its scalability.

To reduce the need for manual annotations, recent works [44, 45] have introduced the use of the Segment Anything Model (SAM) [17] in combination with NeRF or 3DGS. Specifically, OR-NeRF employs Grounded-SAM [32] to locate a single-view 2D inpainting mask for the object to be removed. It then projects 3D points of the object’s surface into other views, which are used as prompts for

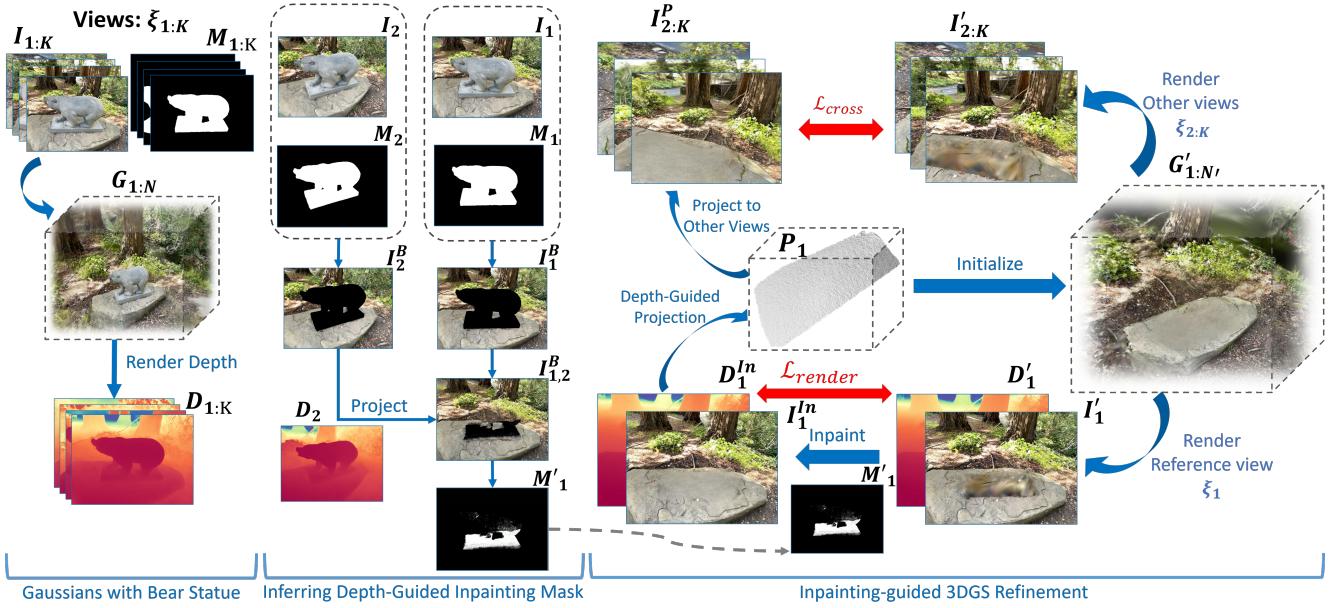


Figure 2. Overview of 3D Gaussian Inpainting with Depth-Guided Cross-View Consistency. Given a 3D Gaussian Splatting model $G_{1:N}$ pretrained on multi-view images $I_{1:K}$ at camera poses $\xi_{1:K}$, our goal is to perform 3D inpainting based on the object masks $M_{1:K}$ (e.g., provided by SAM). With the rendered depth maps $D_{1:K}$, the stage of Inferring Depth-Guide Inpainting Mask is able to refine the inpainting masks to preserve visible backgrounds across camera views. The stage of Inpainting-guided 3DGS Refinement then utilizes such masks to jointly update the new Gaussians $G'_{1:N'}$ for both novel-view rendering and inpainting purposes.

SAM to generate masks for the remaining views. Similarly, Gaussian Grouping [44] enhances 3DGS by incorporating semantic feature learning, allowing the model to jointly render RGB images and segmentation maps, where the segmentation supervision is derived from SAM. While these methods significantly reduce the burden of manual mask creation, they inpaint 2D images of different views separately and optimize the inpainted NeRF by treating all the 2D inpaintings equally. As a result, the above approaches still face difficulties in producing consistent multi-view results, as mentioned in [7, 18, 38]

To alleviate this problem, more advanced approaches [7, 18, 38] focus on improving cross-view consistency. For instance, MALD-NeRF fine-tunes a scene-specific Low-Rank Adaptation (LoRA) [14] module for a pre-trained diffusion model to inpaint images of each scene. By introducing a LoRA module for each scene, the diffusion model can inpaint more consistent content across different views. GScream [38], on the other hand, applies diffusion-based 2D inpainting on a chosen reference view. By predicting the depth map of the inpainted reference view, GScream incorporates cross-view feature consistency between any other view and the reference view, optimizing geometric alignment across views. These methods represent a significant step forward in achieving automatic, consistent 3D inpainting, addressing the practical limitations of earlier approaches. Nonetheless, the aforementioned methods rely

on per-view 2D inpainting masks for 2D inpainting models as input, while some areas in those masks are visible from other views, as noted in [44]. Consequently, the inpainted content for these visible areas may not align with the original scene (as illustrated in the red branch in Figure 1). This inconsistency might be propagated to the inpainted 3D scene, hindering the reliability of their results.

3. Method

3.1. Problem Definition and Model Overview

We begin with the notations and settings of our proposed framework. Given a 3D Gaussian Splatting (3DGS) model [44] $G_{1:N} = \{G_1, G_2, \dots, G_N\}$ (N denotes the number of Gaussians) pretrained for K multi-view images $I_{1:K} = \{I_1, I_2, \dots, I_K\}$ with their camera poses $\xi_{1:K} = \{\xi_1, \xi_2, \dots, \xi_K\}$, our goal is to remove the Gaussians corresponding to a particular object (e.g., the bear statue) described by 2D object masks $M_{1:K} = \{M_1, M_2, \dots, M_K\}$. More precisely, we aim to update the above 3DGS so that the optimized Gaussians $G'_{1:N'}$ (with N' remaining Gaussians) allow novel view rendering without the object of interest presented. Take Figure 2 as an example, the bear statue is to be removed from the scene of interest, and its segmentation masks $M_{1:K}$ from $I_{1:K}$ can be produced by models like SAM [17] (see supplementary materials for details).

To address the above 3D Gaussian Inpainting with Depth-

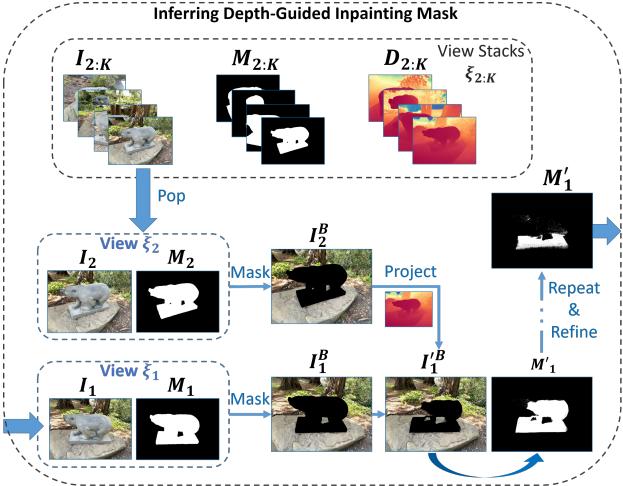


Figure 3. **Inferring Depth-Guided Inpainting Mask.** Taking $\{I_1, M_1\}$ at view ξ_1 as an example reference view, the original background region I_1^B can be first produced. We then project the background region I_2^B from ξ_2 to ξ_1 , updating I_1^B and the associated inpainting mask M'_1 . By repeating this process across camera views, the final inpainting mask M'_1 contains only the regions that are *not* visible at any training camera views.

Guided Cross-View Consistency (3DGIC). Our 3DGIC comprises two learning stages: **Depth-Guided Inpainting Mask** and **Inpainting-guided 3DGS Refinement**. The former refines the object mask guided by both semantics and depth maps observed across $I_{1:K}$, while the latter performs inpainting with cross-view consistency for updating the Gaussians $G'_{1:N'}$.

3.2. Inferring Depth-Guided Inpainting Masks

Given multi-view images $I_{1:K}$ of a scene with binary masks $M_{1:K}$ depicting the object to be removed, we aim to infer a proper mask M' for inpainting images at each view under the guidance of depth images $D_{1:K}$ rendered from $G_{1:N}$. As a result, the masked image $I'^B_{1:K}$ only contains background pixels that are visible at other camera views. The i -th masked image I'^B_i is defined as:

$$I'^B_i = I_i \cdot (\mathbb{1} - M_i), \quad (1)$$

where $\mathbb{1}$ denotes a tensor with the same size as M and all the elements are one.

Take $\{I_1, M_1\}$ in Figure 3 as an example, our process of inferring Depth-Guided Inpainting Masks takes the original image I_1 from ξ_1 and masks out the areas in M_1 as the original visible backgrounds $I_1^B = I_1 \cdot (\mathbb{1} - M_1)$ from view ξ_1 . To explore all the visible background pixels from other views $\xi_{2:K}$, we take $I_{2:K}$ with their masks $M_{2:K}$ and rendered depth $D_{2:K}$ at $\xi_{2:K}$, and we project the above background pixels from each view to ξ_1 . Taking view ξ_2 as an example, the visible backgrounds I_2^B in I_2 ($I_2^B = I_2 \cdot (\mathbb{1} - M_2)$) are

projected into the 3D space via D_2 and ξ_2 and then back-projected to ξ_1 . Among all the back-projected pixels, we consider the pixel coordinates that lie inside M_1 as visible backgrounds from I_2 , denoted as $I_{1,2}^B$. The operation for obtain $I_{1,2}^B$ is calculated as:

$$I_{1,2}^B = Proj^{2D}(Proj^{3D}(I_2^B, D_2, \xi_2), \xi_1) \cdot M_1, \quad (2)$$

where $Proj^{3D}(\cdot, \cdot, \cdot)$ denotes the 3D projection function that projects 2D colored pixels in I_2^B into 3D point clouds with its depth map D_2 and camera pose ξ_2 , while $Proj^{2D}(\cdot, \cdot)$ represents the 2D projection function that projects the 3D colored point cloud back to ξ_1 as colored pixels. With the above operation, the corresponding pixel coordinates of $I_{1,2}^B$ are directly excluded from M_1 , and thus the inpainting mask is refined as M'_1 , and the masked image $I'^B_1 = I_1^B + I_{1,2}^B$ at ξ_1 is further obtained. Similarly, we repeat this process through all the views $\xi_{2:K}$ to infer the final Depth-Guided Inpainting Mask M'_1 and the masked image I'^B_1 at ξ_1 . Also, we produce the depth-guided inpainting masks $M'_{1:K}$ for all the views $\xi_{1:K}$. Please refer to our supplementary material for more details about this inferring process.

It is worth noting that the above process is deterministic. It not only reduces the uncertainty of image regions to inpaint at each view, but it also makes the updating of the 3DGS model for rendering the unpainted scene more effective, as discussed in the following subsection.

3.3. Inpainting-guided 3DGS Refinement

The aim of this stage is to optimize $G'_{1:N'}$ with masked $I'_{1:K}$ obtained at $\xi_{1:K}$ with cross-view consistency so that rendering of the corresponding high-fidelity scene can be produced, realizing the task of 3D inpainting. As shown in Figure 2, the 3DGS for this inpainting scene can be first updated by removing the Gaussians with the semantic labels corresponding to the masked region (e.g., “bear” in the original Gaussian $G_{1:N}$), and replaced by the same amount of randomly initialized Gaussians in the masked region (e.g., with bear removed; see [44] and our supplementary materials).

Take ξ_1 as the reference view for an example, the rendered image I'_1 and depth map D'_1 of $G'_{1:N'}$ at ξ_1 are inpainted by a 2D inpainter [33, 36] (using M'_1 as inpainting mask) as:

$$\begin{aligned} I_1^{In} &= Inpaint_{2D}(I'_1, M'_1) \\ D_1^{In} &= Inpaint_{2D}(D'_1, M'_1), \end{aligned} \quad (3)$$

where $Inpaint_{2D}(\cdot, \cdot)$ denotes the 2D inpainting process, and I_1^{In} and D_1^{In} represents the 2D-inpainted results of I'_1 and D'_1 , respectively. To ensure I'_1 looks identical to I_1^{In} , the *rendering loss* at ξ_1 is defined as:

$$\mathcal{L}_{rendering} = \mathcal{L}_{rgb} + \mathcal{L}_{depth}. \quad (4)$$

Note that the image recovery loss \mathcal{L}_{rgb} is calculated as:

$$\mathcal{L}_{rgb} = \|I'_1 - I_1^{In}\|_1 + \mathcal{L}_{SSIM}(I'_1, I_1^{In}), \quad (5)$$

where the \mathcal{L}_{SSIM} denotes the structure similarity loss [15]. And, the depth loss \mathcal{L}_{depth} is defined as:

$$\mathcal{L}_{depth} = \|D'_1 - D_1^{In}\|_1. \quad (6)$$

To further ensure the masked regions in $I'_{2:K}$ (with respect to $M'_{2:K}$) are cross-view consistent with the 2D-inpainted region in I_1^{In} , we project the inpainted region of I_1^{In} into the 3D space as a set of colored point clouds P_1 , followed by re-projecting back to $\xi_{2:K}$ as supervision. Thus, P_1 is calculated as:

$$P_1 = Proj^{3D}(I_1^{In} \cdot M'_1, D_1^{In}, \xi_1), \quad (7)$$

where $Proj^{3D}(\cdot, \cdot, \cdot)$ is the same projection function in Eqn. 2. For each view ξ_k of $\xi_{2:K}$, the back-projected image I_k^P for supervision is denoted as:

$$I_k^P = I'_k \cdot (1 - M'_k) + Proj^{2D}(P_1, \xi_k) \cdot M'_k, \quad (8)$$

where $Proj^{2D}(\cdot, \cdot)$ is also the same 2D projection function in Eqn. 2. To this end, the cross-view consistent loss \mathcal{L}_{cross} is defined as:

$$\mathcal{L}_{cross} = \sum_{k \in 2 \dots K} \mathcal{L}_{LPIPS}(I'_k, I_k^P), \quad (9)$$

where \mathcal{L}_{LPIPS} denotes the LPIPS [48] loss that calculates the perceptual similarity between I'_k and I_k^P .

Finally, the overall loss for 3D inpainting is calculated by $\mathcal{L}_{inpaint} = \mathcal{L}_{render} + \mathcal{L}_{cross}$. We note that by conducting $\mathcal{L}_{inpaint}$, $G'_{1:N'}$ is guaranteed to inpaint the object-removed 3D scene with cross-view consistency by taking $\{I_1^{In}, D_1^{In}\}$ as guidance.

3.4. Training and Inference

3.4.1. Training

During the training (optimization) process, we calculate the refined mask M' described in Sect. 3.2 for all K views and choose the view with the largest refined mask as the reference view. This is because the 2D inpainted result from this view covers the most 3D space compared to other views, allowing us to provide a more informative cross-view optimization. By choosing the reference view, $\mathcal{L}_{inpaint}$ is applied to optimize $G'_{1:N'}$. To this end, $G'_{1:N'}$ is properly supervised to ensure the 3D scene is reasonably inpainted and consistent across different views.

3.4.2. Inference

Once we finish the optimization of the inpainted scene with our 3DGIC, the optimized Gaussians $G'_{1:N'}$ are able to render a novel view synthesis of the scene by using arbitrary camera poses.

4. Experiments

4.1. Datasets

To evaluate the effectiveness of our method, we conduct experiments on the most used real-world benchmark dataset: the SPIn-NeRF [25] dataset. This dataset contains *ten* real-world scenes, including indoor and outdoor scenes. Each scene is composed of 60 frames of training images and 40 frames of testing images where a certain object in the scene is removed, with camera poses of all 100 images available. The binary mask of the object to be removed is also provided in each frame for evaluation. Following the setting of [7, 18, 25, 38, 44], we resize each image as 1008×567 in resolution for all our experiments and show the comparisons quantitatively and qualitatively.

Since the camera poses in all the scenes provided in the SPIn-NeRF dataset only cover a small range (i.e., all the image frames are captured near the front view of the scene), we additionally include qualitative comparisons with several scenes covering 360° of camera poses to show the effectiveness of our design, specifically for our Depth-Guided In-painting Mask. Following Gaussian Grouping [44], we take the “bear” scene provided in InNeRF360 [37], the “counter” scene in Mip-NeRF360 [2], and the “figureines” scene in LeRF [16] for the additional qualitative evaluations. Since these scenes are not originally for the 3D inpainting task, we manually select an object in each scene as the object to be removed and select the corresponding ID in the segmentation map obtained from SAM [17] as the object mask in each view. Please refer to our supplementary material for a detailed description of these scenes.

4.2. Quantitative Evaluations

Table 1 shows the comparisons between our 3DGIC (with LAMA [36] or LDM [33] as 2D inpainter) and several state-of-the-art approaches such as SPIn-NeRF [25], MVIP-NeRF [7], Gaussian Grouping [44], MALD-NeRF [18], and GScream [38] using the SPIn-NeRF dataset. Following SPIn-NeRF and MALD-NeRF, we conduct FID [13], masked FID (m-FID), LPIPS [48], and masked LPIPS (m-LPIPS) as our evaluation matrices, where m-FID and m-LPIPS calculate the FID and LPIPS scores only inside the ground truth inpainting masks. We note that the official implementation of MALD-NeRF is currently unavailable; we directly use the output results provided on their official project page for evaluation. As for other state-of-the-arts, we reproduce results from their official implementations and the released configurations.

From Table 1, we can see that the LDM version of our 3DGIC achieves the best score on all four evaluation matrices. As for our 3DGIC using a non-diffusion-based model of LAMA as the 2D inpainter, the results still outperform MVIP-NeRF and MALD-NeRF, where both use LDM as the

Table 1. **Quantitative evaluation on the SPIn-NeRF dataset in terms of FID and LPIPS.** Note that m-FID and m-LPIPS represent that the FID and LPIPS scores are only calculated within the ground truth inpainting masks.

	Representation	2D inpainter	FID \downarrow	m-FID \downarrow	LPIPS \downarrow	m-LPIPS \downarrow
SPIn-NeRF [25]	NeRF	LAMA [36]	49.6	153.4	0.31	0.053
MVIP-NeRF [7]	NeRF	LDM [33]	50.5	173.4	0.31	0.050
Gaussian Grouping [44]	Gaussian Splatting	LAMA [36]	44.7	132.5	0.30	0.037
MALD-NeRF [18]	NeRF	LDM [33]	44.9	113.5	0.26	0.031
GScream [38]	Gaussian Splatting	LDM [33]	38.6	101.6	0.28	0.033
3DGIC (Ours)	Gaussian Splatting	LAMA [36]	41.7	102.4	0.28	0.032
3DGIC (Ours)	Gaussian Splatting	LDM [33]	36.4	96.3	0.26	0.028

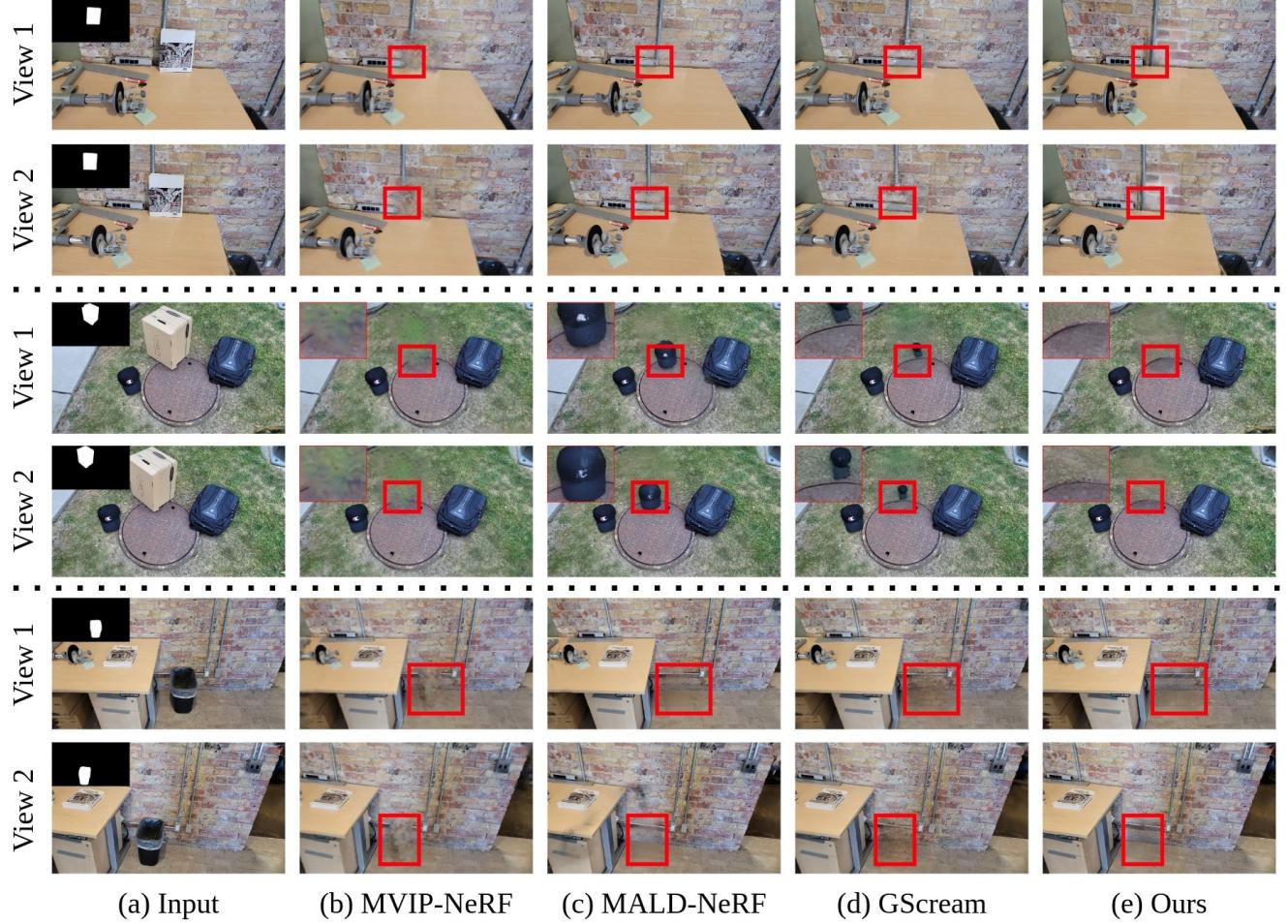


Figure 4. **Qualitative results on the SPIn-NeRF [25] dataset.** Two different views of the same scene are shown for each inpainting example. We compare rendering results against MVIP-NeRF [7], MALD-NeRF [18], and GScream [38]. We can see from the regions highlighted by the red boxes that our 3DGIC performs better in terms of multi-view consistency and rendering fidelity

inpainter. The above results show that while using a better 2D inpainter achieves better results, the improvements in our 3DGIC do not come solely from a better 2D inpainter. This suggests that our model is not bundled by 2D inpainters and achieves 3D inpainting with improved fidelity.

4.3. Qualitative Results

In Figure 4, we qualitatively compare our 3DGIC with MVIP-NeRF [7], MALD-NeRF [18], and GScream [38] using the testing set of SPIn-NeRF dataset. In this figure, each of the two rows shows the results of the same scene with different viewpoints, while the first column shows the images containing the object to be removed along with the object masks at the upper-left corner. Specifically, from the

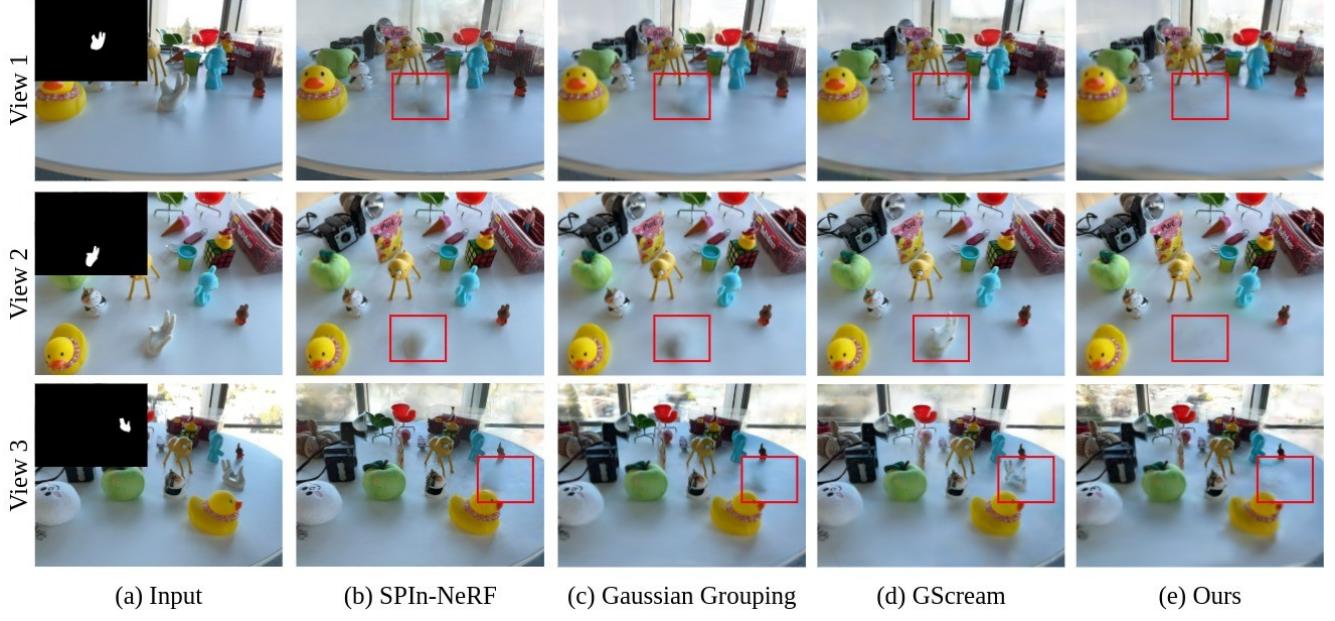


Figure 5. Qualitative results on the *Figurines* scene from the LeRF [16] dataset. We compare the rendering results with SPIn-NeRF [25], Gaussian Grouping [44], and GScream [38]. The three rows show different views of the scene, whereas the first column shows the input images with the object masks of the unwanted object. The regions highlighted by the red boxes show that our 3DGIC inpaints a smoother table surface without artifacts.

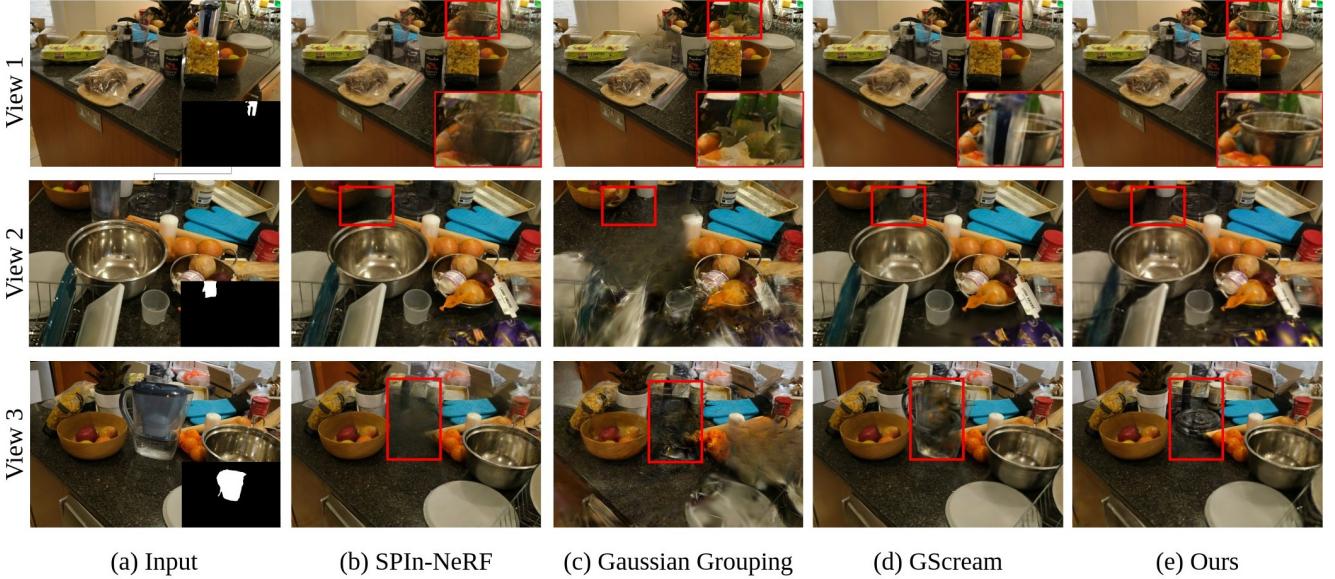


Figure 6. Qualitative results on the *Counter* scene from the MipNeRF360 [2] dataset. We compare the rendering results with SPIn-NeRF [25], Gaussian Grouping [44], and GScream [38]. The three rows show different views of the scene, where we zoom in a certain region in the first row to highlight the difference between each method. We can see from the regions highlighted by the red boxes that our 3DGIC correctly inpaints the water bottle without manipulating any other objects on the table (e.g., the plastic cover).

first two rows, we observe that while GScream and MALD-NeRF both show high-fidelity images, some of the visible details from the input image (e.g., the electrical socket on the table) are not preserved properly. For the third and fourth

rows, where we zoom in on certain areas inside the red boxes, although it is reasonable for MALD-NeRF to generate a hat in the inpainted region, the logo on the hat is not consistent across different views. As for MVIP-NeRF, blurry images

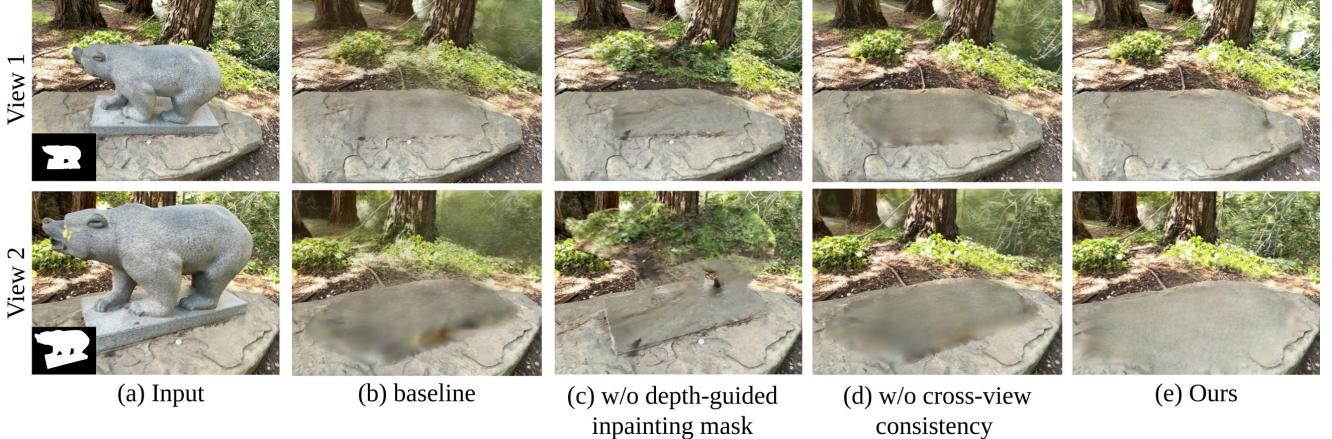


Figure 7. **Ablation studies on the *Bear* scene from the InNeRF360 [37] dataset.** We verify the effectiveness of our Inferring Depth-Guided Inpainting Mask and Inpainting-guided 3DGS Refinement.

are generated in all cases. Oppositely, our 3DGIC generates high-fidelity images with multi-view consistency and preservation of the visible backgrounds.

In Figure 5 and Figure 6, we further show the qualitative comparisons with SPIIn-NeRF, Gaussian Grouping, and GScream using the *Figurines* dataset from LeRF [16] and the *Counter* dataset from MipNeRF360 [2], where each shows results from three different viewpoints. For Figure 5, we can see that both SPIIn-NeRF and Gaussian Grouping leave obvious black holes and shadows in the inpainting region, while GScream does not clearly remove the object of interest. In contrast, our 3DGIC successfully removes the unwanted object and produces smooth and multi-view consistent results without leaving heavy shadows. For Figure 6, where certain areas are cropped by red boxes and zoomed in in the first row, GScream does not fully remove the object of interest either. SPIIn-NeRF not only removes the object of interest but also inpaints other objects in the background. As for Gaussian Grouping, which uses GroundedSAM [32] to detect inpainting mask with the text prompt “*blurry hole*” as input, the GroundedSAM model locates other regions rather than focusing on the object removed region, producing blurry and inconsistent inpainting results across different views. In the contrary, our 3DGIC locates the regions to be inpaint properly and hence produces high-fidelity results while preserving all the other background objects.

4.4. Ablation Study

To further analyze the effectiveness of our designed modules (i.e., Inferring Depth-Guided Masks and Inpainting-guided 3DGS Refinement), we conduct ablation studies on the “*bear*” scene from InNeRF360 [37], as shown in Figure 7. Column (a) shows the input images with the bear statue and their corresponding object mask. The baseline model (b) uses the original object masks as the inpainting

masks and directly applies all the inpainted 2D images as input to fine-tune a 3DGS model. The results of model (b) show blurry contents all over the rendered image, while the inpainted results are not consistent across different views. For model (c), the original object masks are applied as the 2D inpainting model, with our Inpainting-guided 3DGS Refinement. Although the rendered images of model (c) show better fidelity, using the original object masks as inpainting masks results in modifications to the visible backgrounds. For model (d), our inferred depth-guided masks $M'_{1:K}$ are applied as the 2D inpainting masks, but all the 2D inpainting results are directly used as inputs to fine-tune the 3DGS model. As a result, although the backgrounds are preserved, the inpainted region is blurry and not consistent across the views. As for our full model in the last column (e), the depth-guided masks are used, and the 3D Inpainting with Cross-View Consistency is applied, achieving the best results. This verifies the success of our proposed modules and strategies for 3D inpainting.

5. Conclusions

In this paper, we propose the 3D Gaussian Inpainting with Depth-Guided Cross-View Consistency (3DGIC) for inpainting real-world 3D scenes represented by 3D Gaussian Splatting (3DGS) models. With the conduction of our Inferring Depth-Guided Inpainting Masks, we are allowed to obtain precise inpainting masks by considering rendered depth maps and visible background information from other views. With these depth-guided inpainting masks properly obtained, our Inpainting-guided 3DGS Refinement optimizes a newly initialized 3DGS model and performs 3D inpainting simultaneously. In our experiments, we quantitatively and qualitatively show that our 3DGIC is able to handle scenes with various ranges of camera views and perform favorably against existing 3D inpainting approaches.

3D Gaussian Inpainting with Depth-Guided Cross-View Consistency

Supplementary Material

A. Additional Details of 3DGIC

A.1. Details of Backbone 3D Gaussian Splatting Model

Given the multi-view images $I_{1:K}$ with corresponding camera poses $\xi_{1:K}$ of a 3D scene, the vanilla 3DGS [15] model parameterize each Gaussian G_i in $G_{1:N}$ with its 3-dimensional centroid $\mathbf{p}_i \in \mathbb{R}^3$, a 3-dimensional standard deviation $\mathbf{s}_i \in \mathbb{R}^3$, a 4-dimensional rotational quaternion $\mathbf{q}_i \in \mathbb{R}^4$, an opacity $\alpha_i \in [0, 1]$, and color coefficients \mathbf{c}_i for spherical harmonics in degree of 3. Hence, G_i is represented with a set of the above parameters (i.e., $G_i = \{\mathbf{p}_i, \mathbf{s}_i, \mathbf{q}_i, \alpha_i, \mathbf{c}_i\}$). However, to make sure the 3DGS models in this paper are capable of removing Gaussians corresponding to any indicated object (e.g., “bear” in Figure 2) as described in Sect. 3.3, we incorporate the use of a semantic-aware 3DGS (i.e., Gaussian Grouping [44]) approach as the main backbone 3DGS model of our method. Also, since the rendered depth maps $D_{1:K}$ are utilized as important guidance in our 3DGIC, we additionally combine the use of Relightable Gaussian [11], which produces better depth estimations from 3DGS model as our final backbone for Sect. 3. We now briefly discuss both methods.

Incorporating Semantic Segmentation via Gaussian Grouping. To overcome the lack of fine-grained scene understanding in 3DGS, Gaussian Grouping [44] extends 3DGS by incorporating segmentation capabilities. Along with $I_{1:K}$, Gaussian Grouping additionally takes the Segment Anything Model (SAM) to produce 2D semantic segmentation masks $S_{1:K} = \{S_1, S_2, \dots, S_K\}$ from multiple views as inputs, and an additional 16-dimensional parameter $\mathbf{e}_i \in \mathbb{R}^{16}$ is introduced to represent a 3D Identity Encoding for each Gaussian G_i . Therefore, each Gaussian G_i is extended as $G_i = \{\mathbf{p}_i, \mathbf{s}_i, \mathbf{q}_i, \alpha_i, \mathbf{c}_i, \mathbf{e}_i\}$. To make sure $G_{1:K}$ learns to segment each object represented by $S_{1:K}$ in the scene, a 2D identity loss \mathcal{L}_{id} is applied by calculating cross-entropy between $\hat{S}_{1:K}$ and $S_{1:K}$, where $\hat{S}_{1:K} = \{\hat{S}_1, \hat{S}_2, \dots, \hat{S}_K\}$ denotes the rendered segmentation maps from $G_{1:K}$. Additionally, to further ensure that the Gaussians having the same identities are grouped together, a 3D regularization loss \mathcal{L}_{3D} is applied to enforce each G_i ’s k-nearest 3D spatial neighbors to be close in their feature distance of Identity Encodings. Please refer to the original paper [44] for detailed formulations of segmentation map rendering and \mathcal{L}_{3D} . The design of Gaussian Grouping ensures that the segmentation results are coherent across multiple views, enabling the automatic generation of binary masks for any queried object in the scene.

Produce Reliable Depth Estimations with Relightable Gaussians.

Different from Gaussian Grouping, Relightable Gaussians [11] extends the capabilities of Gaussian Splatting by incorporating Disney-BRDF [4] decomposition and ray tracing to achieve realistic point cloud relighting. Unlike traditional Gaussian Splatting, which primarily focuses on appearance and geometry modeling, Relightable Gaussians also aim to model the physical interaction of light with different surfaces in the scene. Specifically, for each Gaussian G_i , the original color coefficients \mathbf{c}_i is decomposed into a 3-dimensional base color $\mathbf{b}_i \in [0, 1]^3$, a 1-dimensional roughness $r \in [0, 1]$, and incident light coefficients \mathbf{l}_i for spherical harmonics in degree of 3. Subsequently, the Physical-Based Rendering (PBR) process and a point-based ray tracing are applied to obtain the colored PBR 2D images $\hat{I}_{1:K}^{PBR}$ and additionally supervised by $I_{1:K}$. Besides the above extensions on PBR for relighting, Relightable Gaussians also introduces a 3-dimensional normal \mathbf{n}_i for G_i and leverages several techniques, including an unsupervised estimation of a depth map D_i from each input view ξ_i , to enhance the geometry accuracy and smoothness. By conducting this self-supervised estimation and regularization of normal maps and depth maps, the predicted depth map D_i is more reliable than the vanilla 3DGS. Please refer to the original paper of Relightable Gaussians [11] for detailed explanations.

In conclusion, each Gaussian of our 3DGIC is parameterized as $G_i = \{\mathbf{p}_i, \mathbf{s}_i, \mathbf{q}_i, \alpha_i, \mathbf{c}_i, \mathbf{e}_i, \mathbf{b}_i, r, \mathbf{l}_i, \mathbf{n}_i\}$. By combining these methods, we are able to perform reliable depth estimations and effective removal of the Gaussians corresponding to any object in the scene for our 3DGIC.

A.2. Additional Details of Inferring Depth-Guided Inpainting Masks

In Sect. 2.2 in our main paper, we introduce infer proper inpainting masks $M'_{1:K}$ to determine the region to be inpaint by realizing visible background regions across different views. In our implementation, after updating the inpainting masks $M'_{1:K}$ with the process described in Sect. 3.2, we further conduct a refinement for each mask as a post-processing to prevent noisy mask. Taking M'_1 as an example, this process updates M'_1 as:

$$M'_1 \leftarrow Open(M'_1), \quad (10)$$

where $Open(\cdot)$ represents a morphological opening process to reduce noises. This refinement process ensures that small noisy pixels are suppressed in our Depth-Guided Inpainting Masks.

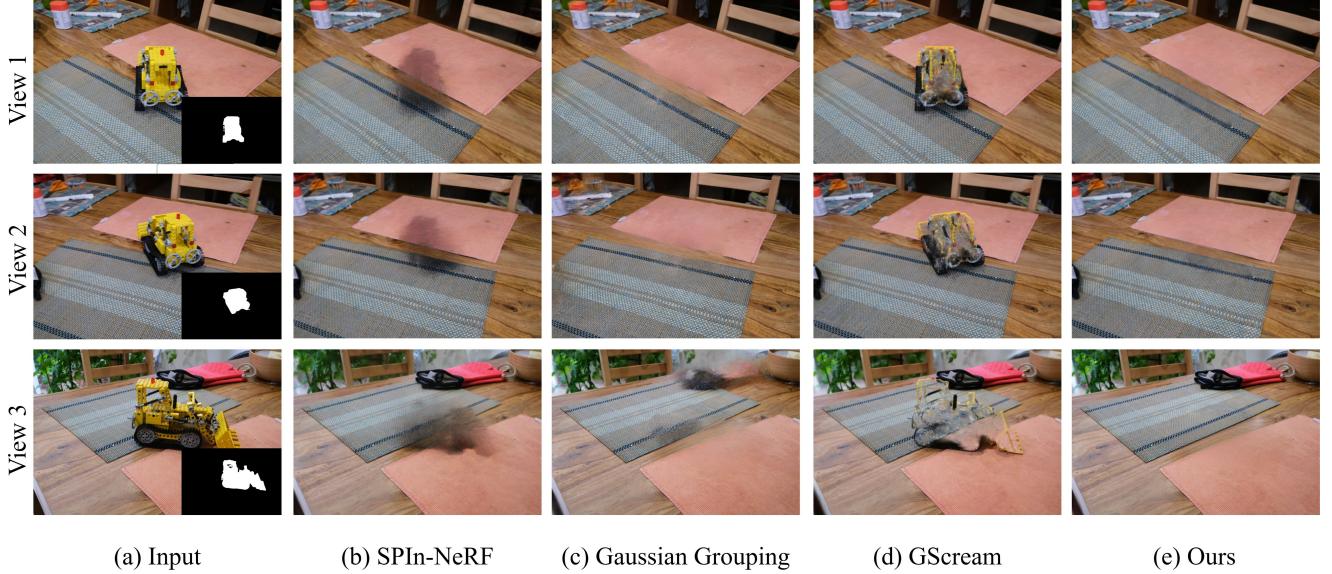


Figure A8. **Qualitative results on the *Kitchen* scene from the MipNeRF360 [2] dataset.** We compare the rendering results with SPIn-NeRF [25], Gaussian Grouping [44], and GScream [38]. The three rows show different views of the scene. We can see that our 3DGIC inpaint a smooth kitchen table, while other approaches produce blurry results.

A.3. Additional Details of Initializing Inpainted Gaussian

In Sect. 3.3, we introduce to remove the Gaussians with semantic labels corresponding to the “bear” object in $G_{1:N}$ and replace by the same amount of randomly initialized Gaussians in the masked region as the initialization of $G'_{1:N'}$. We now detail this initialization process for $G'_{1:N'}$.

When first removing the Gaussians corresponding to the “bear” object, we directly use the remaining Gaussian to render the image I'_1 and depth map D'_1 . Following the 2D inpainting process described in Sect. 3.3, the inpainted image I_1^{In} and depth map D_1^{In} are produced and projected into 3D space as colored point clouds P_1 . We then use the 3D coordinates of P_1 as the initialized 3D position for the newly introduced Gaussians for $G'_{1:N'}$, since P_1 represents the ideal surface of the inpainted 3D Gaussian provided by I_1^{In} after removing the bear. Note that if the number of points in P_1 does not match the number of newly initialized Gaussians in $G'_{1:N'}$ (also the number of removed Gaussians in $G_{1:N}$), we apply random selection to the coordinates of P_1 to match the number of the newly introduced Gaussians. As for the other parameters of the newly introduced Gaussians in $G'_{1:N'}$, we follow Gaussian Grouping [44] to average the parameters of each Gaussian’s 5-nearest neighbors (in 3D space) from the remaining Gaussians as initialization. By this process, $G'_{1:N'}$ is properly initialized.

A.4. Implementation Details

In all our experiments, we train one model for each object category, using a single NVIDIA RTX 3090 GPU (24G) for

training with the PyTorch [28] libraries. For each scene, 5000 iterations of optimization are applied to obtain the inpainted 3DGS model. We also use the official implementation of [7, 25, 38, 44] for comparison. When applying 2D inpainting models to the image and depth map to be inpaint, if we use non-diffusion-based LAMA [36] as inpainter, the RGB image and depth map are inpainted separately. However, if LDM [33] is applied as our 2D inpainters, we follow the suggestion in NeRFiller [39] to stack the RGB image and the depth map in the same image for inpainting to ensure the inpainted RGB image and the depth map are consistent in terms of the geometry details. Specifically, we crop a 512×512 patch for the RGB image and the depth map to be inpainted center at the pixel coordinate of the inpainting mask’s center, and paste the cropped RGB patch to a 1024×1024 -resolution black image at the upper right corner with the cropped depth map at the lower left corner as the input image for the LDM. Similarly, we also crop a 512×512 patch for the inpaint masks and put them to the upper right and lower left corner of another 1024×1024 -resolution black image as the input binary inpainting mask for the LDM. We then use the prompt “*an RGB image and a depth image of the same scene*” to inpaint the input image. Finally, the inpainted RGB patch and the depth map patch are pasted back to the original image and depth map, respectively, as the 2D inpainting result. It is worth noting that we apply the 2D inpainting process for every 500 iterations. Following MALD-NeRF [18], we use the technique of partial DDIM [34], to start from latter step of the denoising process as optimization iteration grows. Specifically, for a 50-step

DDIM process, we start from step 0 of the LDM denoising process for step 0 of our optimization. After 500 iteration steps, the second time of the LDM inpainting starts from step 5 of the DDIM process and so on. When our optimization reaches the last 500 iterations, the 2D inpainting process only denoises using the last five steps of DDIM. This prevents inpainting results that are too different from the current scene and provides more stability for our optimization process.

A.5. Dataset Details

For the “figurines” scene from LeRF [16] dataset, we have 260 training frames and 40 testing frames, each with a resolution of 986×728 . For the “bear” dataset from In-NeRF360 [37], we have 90 training frames and 6 testing frames, each with a resolution of 985×729 . As for “counter” and “kitchen” scenes from MipNeRF360 [2], 240 (230 for training and 10 for testing) and 279 (270 for training and 9 for testing) frames are available in total, respectively. Both scenes are in the resolution of 779.

B. Additional Experiments

We additionally show the results on the “kitchen” scene from the MipNeRF360 [2] dataset in Figure A8. We can see that SPIIn-NeRF produces blurry result, while GScream fail to handle camera views with a wide range and not able to remove the excavator clearly. Although Gaussian Grouping also produces plausible results at the excavator-removed regions, it incorrectly detects the glove behind the excavator as region to be inpaint by using the “black blurry hole” as the prompt for Gounded-SAM [32] to find inpainting masks and therefore changes the background that should not be changed (shown in the third view). On the other hand, our 3DGIC locates the proper region to inpaint and produces smooth and high fidelity results.

C. Limitations

We now discuss the potential limitations of our 3DGIC. Since our 3DGIC uses the rendered depth map as guidance for the 3D inpainting process, the reliability of the rendered depth map becomes an important issue. As detailed in Sect. A.1, we combine the optimization technique introduced in Relightable Gaussians [11] to conduct a self-supervised loss for the predicted normal map and the rendered depth map to enhance the accuracy of the rendered depth map. However, if the input views are too sparse, the rendered depth map would not be guaranteed to be accurate, which hinders the inferring of Depth-Guided Mask and the achievement of cross-view consistency. Another potential limitation of our 3DGIC lies in the capability of the SAM [17] model. As detailed in Sect. A.1, we use SAM to produce 2D segmentation masks and use these masks as supervision for our backbone 3DGS model so that we don’t have to manually annotate the 2D object mask of the object to be removed like SPIIn-NeRF [25]. However, if the object

to be removed is too small, the SAM model would confuse it with other objects and not produce the correct segmentation mask for the object. To overcome the above limitations, studies on the production of reliable depth maps for 3DGS models with only sparse input views and producing a more accurate segmentation mask for any object would be possible directions to improve the quality of 3D Gaussian inpainting.

References

- [1] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. 1
- [2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 5, 7, 8, 2, 3
- [3] Wolfgang Broll. Augmented reality. In *Virtual and Augmented Reality (VR/AR) Foundations and Methods of Extended Realities (XR)*, pages 291–329. Springer, 2022. 1
- [4] Brent Burley and Walt Disney Animation Studios. Physically-based shading at disney. In *AcM Siggraph*, 2012. 1
- [5] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022. 1
- [6] Guikun Chen and Wenguan Wang. A survey on 3d gaussian splatting. *arXiv preprint arXiv:2401.03890*, 2024. 1
- [7] Honghua Chen, Chen Change Loy, and Xingang Pan. Mvip-nerf: Multi-view 3d inpainting on nerf scenes via diffusion prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 3, 5, 6
- [8] Ciprian Corneanu, Raghudeep Gadde, and Aleix M Martinez. Latentpaint: Image inpainting in latent space with diffusion models. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2024. 1
- [9] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022. 2
- [10] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. 1, 2
- [11] Jian Gao, Chun Gu, Youtian Lin, Hao Zhu, Xun Cao, Li Zhang, and Yao Yao. Relightable 3d gaussian: Real-time point cloud relighting with brdf decomposition and ray tracing. *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 1, 3
- [12] Yuying Hao, Yi Liu, Zewu Wu, Lin Han, Yizhou Chen, Guowei Chen, Lutao Chu, Shiyu Tang, Zhiliang Yu, Zeyu

- Chen, et al. Edgeflow: Achieving practical interactive segmentation with edge-guided flow. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 1
- [13] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 5
- [14] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 3
- [15] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)*, 2023. 1, 2, 5
- [16] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2023. 5, 7, 8, 3
- [17] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2023. 2, 3, 5
- [18] Chieh Hubert Lin, Changil Kim, Jia-Bin Huang, Qinbo Li, Chih-Yao Ma, Johannes Kopf, Ming-Hsuan Yang, and Hung-Yu Tseng. Taming latent diffusion model for neural radiance field inpainting. *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 1, 2, 3, 5, 6
- [19] Zhiheng Liu, Hao Ouyang, Qiyu Wang, Ka Leong Cheng, Jie Xiao, Kai Zhu, Nan Xue, Yu Liu, Yujun Shen, and Yang Cao. Infusion: Inpainting 3d gaussians via learning depth completion from diffusion prior. *arXiv preprint arXiv:2404.11613*, 2024. 2
- [20] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1
- [21] Márcio CF Macedo and Antonio L Apolinario. Occlusion handling in augmented reality: past, present and future. *IEEE Transactions on Visualization and Computer Graphics*, 29(2): 1590–1609, 2021. 1
- [22] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 1
- [23] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2
- [24] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Marcus A Brubaker, Jonathan Kelly, Alex Levinstein, Konstantinos G Derpanis, and Igor Gilitschenski. Reference-guided controllable inpainting of neural radiance fields. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2023. 1
- [25] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Konstantinos G Derpanis, Jonathan Kelly, Marcus A Brubaker, Igor Gilitschenski, and Alex Levinstein. Spin-nerf: Multiview segmentation and perceptual inpainting with neural radiance fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 2, 5, 6, 7, 3
- [26] Ashkan Mirzaei, Riccardo De Lutio, Seung Wook Kim, David Acuna, Jonathan Kelly, Sanja Fidler, Igor Gilitschenski, and Zan Gojcic. Reffusion: Reference adapted diffusion models for 3d scene inpainting. *arXiv preprint arXiv:2404.10765*, 2024. 2
- [27] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 1, 2
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 2
- [29] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sndl: Improving latent diffusion models for high-resolution image synthesis. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024. 1
- [30] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d language gaussian splatting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 1
- [31] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14335–14345, 2021. 1
- [32] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159*, 2024. 2, 8, 3
- [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 4, 5, 6, 2
- [34] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. 2
- [35] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022. 1, 2

- [36] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lemitsky. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2022. 1, 4, 5, 6, 2
- [37] Dongqing Wang, Tong Zhang, Alaa Abboud, and Sabine Süsstrunk. Innerf360: Text-guided 3d-consistent object inpainting on 360-degree neural radiance fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 5, 8, 3
- [38] Yuxin Wang, Qianyi Wu, Guofeng Zhang, and Dan Xu. Gscream: Learning 3d geometry and feature consistent gaussian splatting for object removal. *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 1, 2, 3, 5, 6, 7
- [39] Ethan Weber, Aleksander Holynski, Varun Jampani, Saurabh Saxena, Noah Snavely, Abhishek Kar, and Angjoo Kanazawa. Nerfiller: Completing scenes via generative 3d inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20731–20741, 2024. 2
- [40] Silvan Weder, Guillermo Garcia-Hernando, Aron Monszpart, Marc Pollefeys, Gabriel J Brostow, Michael Firman, and Sara Vicente. Removing objects from neural radiance fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1
- [41] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1
- [42] Binbin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. Paint by example: Exemplar-based image editing with diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [43] Shiyuan Yang, Xiaodong Chen, and Jing Liao. Uni-paint: A unified framework for multimodal image inpainting with pretrained diffusion model. In *Proceedings of the 31st ACM International Conference on Multimedia*, 2023. 1
- [44] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3d scenes. *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 1, 2, 3, 4, 5, 6, 7
- [45] Youtan Yin, Zhoujie Fu, Fan Yang, and Guosheng Lin. Or-nerf: Object removing from 3d scenes guided by multiview segmentation with neural radiance fields. *arXiv preprint arXiv:2305.10503*, 2023. 1, 2
- [46] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. 1, 2
- [47] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d gaussian splatting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 1
- [48] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 5