

# ReplaceAnything3D: Text-Guided 3D Scene Editing with Compositional Neural Radiance Fields

Edward Bartrum<sup>1,2\*</sup>Thu Nguyen-Phuoc<sup>3</sup>Chris Xie<sup>3</sup>Zhengqin Li<sup>3</sup>Numair Khan<sup>3</sup>Armen Avetisyan<sup>3</sup>Douglas Lanman<sup>3</sup>Lei Xiao<sup>3</sup><sup>1</sup> University College London <sup>2</sup> Alan Turing Institute <sup>3</sup> Reality Labs Research, Meta

Figure 1. Our method enables prompt-driven object replacement for a variety of realistic 3D scenes.

## Abstract

We introduce *ReplaceAnything3D* model (RAM3D), a novel text-guided 3D scene editing method that enables the

replacement of specific objects within a scene. Given multi-view images of a scene, a text prompt describing the object to replace, and a text prompt describing the new object, our Erase-and-Replace approach can effectively swap objects in the scene with newly generated content while maintaining 3D consistency across multiple viewpoints. We demon-

\*Work done during internship at Meta Reality Labs Research  
Project page: <https://replaceanything3d.github.io>

*strate the versatility of ReplaceAnything3D by applying it to various realistic 3D scenes, showcasing results of modified foreground objects that are well-integrated with the rest of the scene without affecting its overall integrity.*

## 1. Introduction

The explosion of new social media platforms and display devices has sparked a surge in demand for high-quality 3D content. From immersive games and movies to cutting-edge virtual reality (VR) and mixed reality (MR) applications, there is an increasing need for efficient tools for creating and editing 3D content. While there has been significant progress in 3D reconstruction and generation, 3D editing remain a less-studied area. In this work, we focus on 3D scene manipulation by replacing current objects in the scene with new contents with only natural language prompts from a user. Imagine putting on a VR headset and trying to re-model one’s living room. One can swap out the current sofa with a sleek new design, add some lush greenery, or remove clutter to create a more spacious feel.

In this project, we introduce the ReplaceAnything3D model (RAM3D), a text-guided Erase-and-Replace method for scene editing. RAM3D takes multiview images of a static scene as input, along with text prompts specifying which object to erase and what should replace it. Our approach comprises four key steps: 1) we use LangSAM [24] with the text prompts to detect and segment the object to be erased. 2) To erase the object, we propose a text-guided 3D inpainting technique to fill in the background region obscured by the removed object. 3) Next, a similar text-guided 3D inpainting technique is used to generate a new object(s) that matches the input text description. Importantly, this is done such that the mass of the object is minimal. 4) Finally, the newly generated object is seamlessly composited onto the inpainted background in training views to obtain consistent multiview images of an edited 3D scene. Then a NeRF [26] can be trained on these new multiview images to obtain a 3D representation of the edited scene for novel view synthesis. We show that this compositional structure greatly improves the visual quality of both the background and foreground in the edited scene.

Compared to 2D images, replacing objects in 3D scenes is much more challenging due to the requirement for multi-view consistency. Naively applying 2D methods for masking and inpainting leads to incoherent results due to visual inconsistencies in each inpainted viewpoint. To address this challenge, we propose combining the prior knowledge of large-scale image diffusion models, specifically a text-guided image inpainting model, with learned 3D scene representations. To generate new multi-view consistent 3D objects, we adapt Hifa [57], a text-to-3D distillation approach, to our 3D inpainting framework. Compared to

pure text-to-3D approaches, ReplaceAnything3D needs to generate new contents that not only follow the input text prompt but also are compatible with the appearance of the rest of the scene. By combining a pre-trained text-guided image inpainting model with a compositional scene structure, ReplaceAnything3D can generate coherent edited 3D scenes with new objects seamlessly blended with the rest of the original scene.

In summary, our contributions are:

- We introduce an Erase-and-Replace approach to 3D scene editing that enables the replacement of specific objects within a scene at high-resolutions.
- We propose a multi-stage approach that enables not only object replacement but also removal and multiple object additions.
- We demonstrate that ReplaceAnything3D can generate 3D consistent results on multiple scene types, including forward-facing and 360° scenes.

## 2. Related work

**Diffusion model for text-guided image editing** Diffusion models trained on extensive text-image datasets have demonstrated remarkable results, showcasing their ability to capture intricate semantics from text prompts [38, 40, 42]. As a result, these models provide strong priors for various text-guided image editing tasks [6, 11, 18, 30, 31]. In particular, methods for text-guided image inpainting [1, 2] enable local image editing by replacing masked regions with new content that seamlessly blends with the rest of the image, allowing for object removal, replacement, and addition. These methods are direct 2D counterparts to our approach for 3D scenes, where each view can be treated as an image inpainting task. However, 3D scenes present additional challenges, such as the requirement for multi-view consistency and memory constraints due to the underlying 3D representations. In this work, ReplaceAnything3D addresses these challenges by combining a pre-trained image inpainting model with compositional 3D representations.

**Neural radiance fields editing** Recent advancements in NeRFs have led to significant improvements in visual quality [3, 5, 17], training and inference speed [9, 14, 32], and robustness to noisy or sparse input [19, 33, 53, 55]. However, editing NeRFs remains a challenging area. Most of the existing work focuses on editing objects’ appearance or geometry [13, 22, 45, 48, 56]. For scene-level editing, recent works primarily address object removal tasks for forward-facing scenes [27, 29, 52]. Instruct-NeRF2NeRF [10] offers a comprehensive approach to both appearance editing and object addition. However, it modifies the entire scene, while Blended-Nerf [45] and DreamEditor [58] allow for localized object editing but do not support object removal. The work closest to ours is by Mirzaei et al. [27], which can



remove and replace objects using one single image reference from the user. However, since this method relies only on one inpainted image, it cannot handle regions with large occlusions across different views, and thus is only applied on forward-facing scenes.

It is important to note that ReplaceAnything3D adopts an Erase-and-Replace approach for localized scene editing, instead of modifying the existing geometry or appearance of the scene’s contents. This makes ReplaceAnything3D the first method that holistically offers localized object removal, replacement, and addition within the same framework.

**Text-to-3D synthesis** With the remarkable success of text-to-image diffusion models, text-to-3D synthesis has garnered increasing attention. Most work in this area focuses on distilling pre-trained text-to-image models into 3D models, starting with the seminal works Dreamfusion [34] and Score Jacobian Chaining (SJC) [49]. Subsequent research has explored various methods to enhance the quality of synthesized objects [20, 25, 51, 57] and disentangle geometry and appearance [7]. Instead of relying solely on pre-trained text-to-image models, recent work has utilized large-scale 3D datasets such as Objaverse [8] to improve the quality of 3D synthesis from text or single images [21, 36].

In this work, we move beyond text-to-3D synthesis by incorporating both text prompts and the surrounding scene information as inputs. This approach introduces additional complexities, such as ensuring the appearance of the 3D object harmoniously blends with the rest of the scene and accurately modeling object-object interactions like occlusion and shadows. By combining HiFA [57], a text-to-3D distillation approach, with a pre-trained text-to-image inpainting model, ReplaceAnything3D aims to create more realistic and coherent 3D scenes that seamlessly integrate the synthesized 3D objects.

### 3. Preliminary

**NeRF** Neural Radiance Fields (NeRFs) [26] is a compact and powerful implicit representation for 3D scene reconstruction and rendering. In particular, NeRF is continuous 5D function whose input is a 3D location  $\mathbf{x}$  and 2D viewing direction  $\mathbf{d}$ , and whose output is an emitted color  $\mathbf{c} = (r, g, b)$  and volume density  $\sigma$ . This function is approximated by a multi-layer perceptron (MLP):  $F_{\Theta} : (\mathbf{x}, \mathbf{d}) \mapsto (\mathbf{c}, \sigma)$ , which is trained using an image-reconstruction loss. To render a pixel, the color and density of multiple points along a camera ray sampled from  $t=0$  to  $D$  are queried from the MLP. These values are accumulated to calculate the final pixel color using volume rendering:

$$\mathbf{C} = \int_0^D \mathcal{T}(t) \cdot \sigma(t) \cdot \mathbf{c}(t) dt \quad (1)$$

During training, a random batch of rays sampled from various viewpoints is used to ensure that the 3D positions of the reconstructed objects are well-constrained. To render a new viewpoint from the optimized scene MLP, a set of rays corresponding to all the pixels in the novel image are sampled and the resulting color values are arranged into a 2D frame.

In this work, we utilize Instant-NGP [32], a more efficient and faster version of NeRF due to its multi-resolution hash encoding. This allows us to handle images with higher resolution and query a larger number of samples along the rendering ray for improved image quality.

**Distilling text-to-image diffusion models** Dreamfusion [34] proposes a technique called score distillation sampling to compute gradients from a 2D pre-trained text-to-image diffusion model, to optimize the parameters of 3D neural radiance fields (NeRF). Recently, HiFA [57] propose an alternative loss formulation, which can be computed explicitly for a Latent Diffusion Model (LDM). Let  $\theta_{scene}$  be the parameters of an implicit 3D scene,  $y$  is a text prompt,  $\epsilon_{\phi}(\mathbf{z}_t, t, y)$  be the pre-trained LDM model with encoder  $E$  and decoder  $D$ ,  $\theta_{scene}$  can be optimized using:

$$\mathcal{L}_{\text{HiFA}}(\phi, \mathbf{z}, \mathbf{x}) = \mathbb{E}_{t, \epsilon} w(t) [\|\mathbf{z} - \hat{\mathbf{z}}\|^2 + \lambda_{RGB} \|\mathbf{x} - \hat{\mathbf{x}}\|^2] \quad (2)$$

where  $\mathbf{z} = E(\mathbf{x})$  is the latent vector by encoding a rendered image  $\mathbf{x}$  of  $\theta_{scene}$  from a camera viewpoint from the training dataset,  $\hat{\mathbf{z}}$  is the estimate of latent vector  $\mathbf{z}$  by the denoiser  $\epsilon_{\phi}$ , and  $\hat{\mathbf{x}} = D(\hat{\mathbf{z}})$  is a recovered image obtained through the decoder  $D$  of the LDM. Note that for brevity, we incorporate coefficients related to timesteps  $t$  to  $w(t)$ .

Here we deviate from the text-to-3D synthesis task where the generated object is solely conditioned on a text prompt. Instead, we consider a collection of scene views as additional inputs for the synthesized object. To achieve this, we utilize HiFA in conjunction with a state-of-the-art text-to-image *inpainting* LDM that has been fine-tuned to generate seamless inpainting regions within an image. This LDM  $\epsilon_{\psi}(\mathbf{z}_t, t, y, \mathbf{m})$  requires not only a text prompt  $y$ , but also a binary mask  $\mathbf{m}$  indicating the area to be filled in.

## 4. Method

### 4.1. Overview

Our training dataset consists of a collection of  $n$  images  $I_i$ , corresponding camera viewpoints  $\mathbf{v}_i$  and a text prompt  $y_{\text{erase}}$  describing the object the user wishes to replace. Using this text prompt we can obtain masks  $\mathbf{m}_i$  corresponding to every image and camera viewpoint. We additionally have a text prompt  $y_{\text{replace}}$  describing a new object to replace the old object. Our goal is to modify the masked object in every image in the dataset to match the text prompt  $y_{\text{replace}}$ , in a multi-view-consistent manner. We can then train any

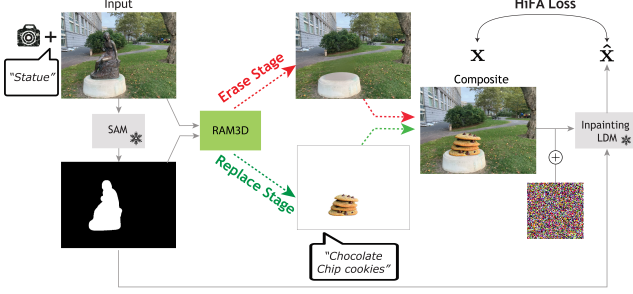


Figure 2. An overview of RAM3D **Erase** and **Replace** stages.

NeRF-like scene reconstruction model using the modified images in order to obtain renderings of the edited scene from novel viewpoints.

Figure 2 illustrates the overall pipeline of our Erase and Replace framework. Instead of modifying existing objects’ geometry and appearance that matches the target text descriptions like other methods [10, 58], we adopt an Erase-and-Replace approach. Firstly, for the **Erase** stage, we remove the masked objects completely and inpaint the occluded region in the background. Secondly, for the **Replace** stage, we generate new objects and composite them to the inpainted background scene, such that the new object blends in with the rest of the background. Finally, we create a new training set using the edited images and camera poses from the original scene, and train a new NeRF for the modified scene for novel view synthesis.

To enable text-guided scene editing, we distill a pre-trained text-to-image inpainting Latent Diffusion Model (LDM) to generate new 3D objects in the scene using HiFA [57]. To address the memory constraints of implicit 3D scenes representations like NeRF, we propose a Bubble-NeRF representation (see Figure 3 and 4) that only models the localised part of the scene that is affected by the editing operation, instead of the whole scene.

#### 4.2. Erase stage

In the Erase stage, we aim to remove the object described by  $y_{erase}$  from the scene and inpaint the occluded background region in a multi-view consistent manner. To do so, we optimise RAM3D parameters  $\theta_{bg}$  which implicitly represent the inpainted background scene. Note that the Erase stage only needs to be performed once for the desired object to remove, after which the Replace stage (Section 4.3) can be used to generate objects or even add new objects to the scene, as demonstrated in the Results section. As a pre-processing step, we use LangSAM [24] with text prompt  $y_{erase}$  to obtain a mask  $\mathbf{m}_i$  for each image in the dataset. We then dilate each  $\mathbf{m}_i$  to obtain *halo* regions  $\mathbf{h}_i$  around the original input mask (see Figure 3).

At each training step, we sample image  $I_i$ , camera  $\mathbf{v}_i$ , mask  $\mathbf{m}_i$ , and halo region  $\mathbf{h}_i$  for a random  $i \in \{1..n\}$ , pro-

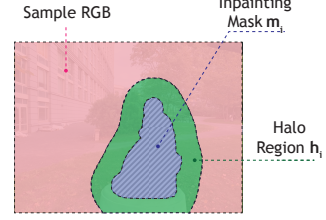


Figure 3. The masked region (blue) serves as a conditioning signal for the LDM, indicating the area to be inpainted. The nearby pixels surrounding  $\mathbf{m}$  form the halo region  $\mathbf{h}$  (green), which is also rendered volumetrically by RAM3D during the Erase stage. The union of these 2 regions is the *Bubble-NeRF* region, whilst the remaining pixels are sampled from the input image (red).

viding them as inputs to RAM3D to compute training losses (left side of Figure 2) (we henceforth drop the subscript  $i$  for clarity). RAM3D volume renders the implicit 3D representation  $\theta_{bg}$  over rays emitted from camera viewpoint  $\mathbf{v}$  which pass through the visible pixels in  $\mathbf{m}$  and  $\mathbf{h}$  (the Bubble-NeRF region). The RGB values of the remaining pixels on the exterior of the Bubble-NeRF are sampled from  $I$  (see Figure 3). These rendered and sampled pixel rgb-values are arranged into a 2D array, and form RAM3D’s inpainting result for the given view,  $\mathbf{x}^{bg}$ . Following the HiFA formulation (see Section 3), we use the frozen LDM’s  $E$  to encode  $\mathbf{x}^{bg}$  to obtain  $\mathbf{z}^{bg}$ , add noise, denoise with  $\epsilon_\psi$  to obtain  $\hat{\mathbf{z}}^{bg}$ , and decode with  $D$  to obtain  $\hat{\mathbf{x}}^{bg}$ . We condition  $\epsilon_\psi$  with  $I$ ,  $\mathbf{m}$  and the empty prompt, since we do not aim to inpaint new content at this stage.

We now use these inputs to compute  $\mathcal{L}_{HiFA}$  (see Equation 2). We next compute  $\mathcal{L}_{recon}$  and  $\mathcal{L}_{vgg}$  on  $\mathbf{h}$  (see Figure 3), guiding the distilled neural field  $\theta_{bg}$  towards an accurate reconstruction of the background.

$$\mathcal{L}_{recon} = MSE(\mathbf{x}^{bg} \odot \mathbf{h}, I \odot \mathbf{h}) \quad (3)$$

$$\mathcal{L}_{vgg} = MSE(vgg_{16}(\mathbf{x}^{bg} \odot \mathbf{h}), vgg_{16}(I \odot \mathbf{h})) \quad (4)$$

This step is critical to ensuring that RAM3D inpaints the background correctly (as shown in Figure 12). Adopting the same formulation as [47], we compute depth regularisation  $\mathcal{L}_{depth}$ , leveraging the geometric prior from a pre-trained depth estimator [39]. In summary, the total loss during the Erase stage is:

$$\mathcal{L}_{Erase} = \mathcal{L}_{HiFA} + \lambda_{recon} \mathcal{L}_{recon} + \lambda_{vgg} \mathcal{L}_{vgg} + \lambda_{depth} \mathcal{L}_{depth} \quad (5)$$

#### 4.3. Replace stage

In the second stage, we aim to add the new object described by  $y_{replace}$  into the inpainted scene. To do so, we optimise the foreground neural field  $\theta_{fg}$  to render  $\mathbf{x}^{fg}$ , which is then composited with  $\mathbf{x}^{bg}$  to form  $\mathbf{x}$ . Unlike  $\theta_{bg}$  in the

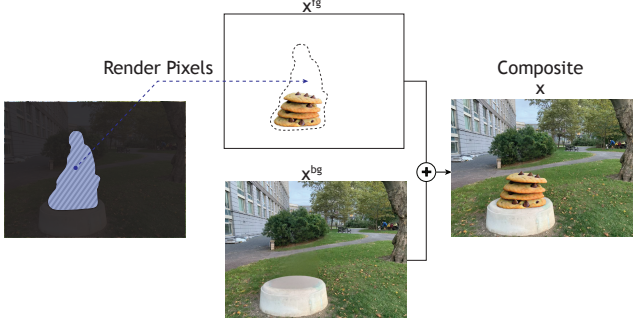


Figure 4. Replace stage: RAM3D volumetrically renders the masked pixels (shown in blue) to give  $\mathbf{x}^{fg}$ . The result is composited with  $\mathbf{x}^{bg}$  to form the combined image  $\mathbf{x}$ .

Erase stage,  $\theta_{fg}$  does not seek to reconstruct the background scene, but instead only the LDM-inpainted content which is located on the interior of  $\mathbf{m}$ . Therefore in the Replace stage, RAM3D does not consider the halo rays which intersect  $\mathbf{h}$ , but only those intersecting  $\mathbf{m}$  (Figure 4). These rendered pixels are arranged in the masked region into a 2D array to give the foreground image  $\mathbf{x}^{fg}$ , whilst the unmasked pixels are assigned an RGB value of 0. The accumulated densities are similarly arranged into a foreground alpha map  $A$ , whilst the unmasked pixels are assigned an alpha value of 0. We now composite the foreground  $\mathbf{x}^{fg}$  with the background  $\mathbf{x}^{bg}$  using alpha blending:

$$\mathbf{x} = A \odot \mathbf{x}^{fg} + (1 - A) \odot \mathbf{x}^{bg} \quad (6)$$

Using the composited result  $\mathbf{x}$ , we compute  $\mathcal{L}_{HiFA}$  as before, but now condition  $\epsilon_\psi$  with the prompt  $y_{replace}$ , which specifies the new object for inpainting. As we no longer require the other losses, we set  $\lambda_{recon}, \lambda_{vgg}, \lambda_{depth}$  to 0.

Since the Erase stage already provides us with a good background, in this stage,  $\theta_{fg}$  only needs to represent the foreground object. To encourage foreground/background disentanglement, on every  $k$ -th training step, we substitute  $\mathbf{x}^{bg}$  with a constant-value RGB tensor, with randomly sampled RGB intensity. This guides the distillation of  $\theta_{fg}$  to only include density for the new object; a critical augmentation to avoid spurious floaters appearing over the background (see Figure 11).

#### 4.4. Training the final NeRF

Once the inpainted background and objects have been generated, we can create a new multi-view dataset by compositing the newly generated object(s) and the inpainted background region for all training viewpoints. We then train a new NeRF, using any variant and framework of choice, to create a 3D representation of the edited scene that can be used for novel view synthesis.

## 5. Results

We conduct experiments on real 3D scenes varying in complexity, ranging from forward-facing scenes to 360° scenes. For forward-facing scenes, we show results for the STATUE and RED-NET scene from SPIn-NeRF dataset [29], as well as the FERN scene from NeRF [26]. For 360° scene, we show results from the GARDEN scene from Mip-NeRF 360°[4]. On each dataset, we train RAM3D with a variety of  $y_{replace}$ , generating a diverse set of edited 3D scenes. Please refer to the [project page](#) for more qualitative results.

**Training details** Each dataset is downsampled to have a shortest image side-length (height) equal to 512, so that square crops provided to the LDM inpainter include the full height of the input image. The FERN scene is an exception, in which we sample a smaller 512 image crop within dataset images with a downsample factor of 2. Details on the resolution and cropping of input images, as well as other implementation details are included in appendices B.4 and B.6.

### 5.1. Qualitative Comparisons

Figures 5, 6 and 7 show qualitative comparison for object replacement by Instruct-NeRF2NeRF [10], Blended-NeRF [2] and the work by Mirzaei et al. [28] respectively.

As shown in Figure 5, Instruct-NeRF2NeRF struggles in cases where the new object is significantly different from the original object (for example, replace the centerpiece with a pineapple or a chess piece in Figure 5 second and third column). More importantly, Instruct-NeRF2NeRF significantly changes the global structure of the scene even when the edit is supposed to be local (for example, only replace the centerpiece with the pineapple). Finally, note that our method is capable of removing objects from the scene completely, while Instruct-NeRF2NeRF cannot (Figure 5 first column).

Figure 6 shows qualitative comparisons with Blended-NeRF. Our method generates much more realistic and detailed objects that blend in much better with the rest of the scene. Meanwhile, Blended-NeRF only focuses on synthesizing completely new objects without taking the surrounding scenes into consideration. The synthesized object therefore looks saturated and outlandish from the rest of the scene. Moreover, due to the memory constraint of CLIP [37] and NeRF, Blended-NeRF only works with images 4-time smaller than ours (2016×1512 vs. 504×378).

Since Mirzaei et al. [28] did not share their code publicly, we report the images adapted from their paper in Figure 7. Our method achieves comparable object replacement results while handling more complex lighting effects such as shadows between the foreground and background objects.





Figure 5. Comparison with Instruct-NeRF2NeRF.

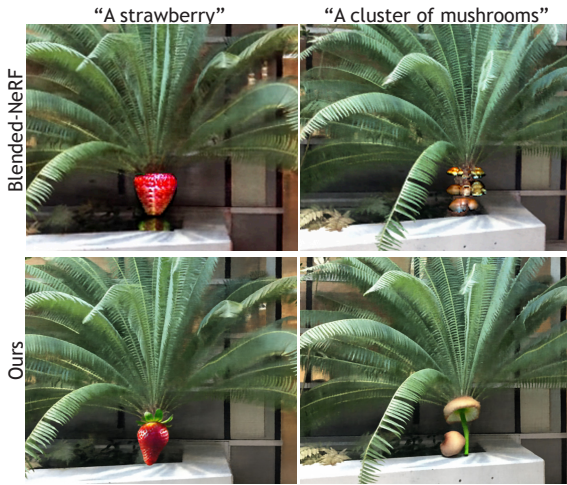


Figure 6. Qualitative comparison with Blended-NeRF for object replacement. Our method generates results with higher quality and capture more realistic lighting and details.

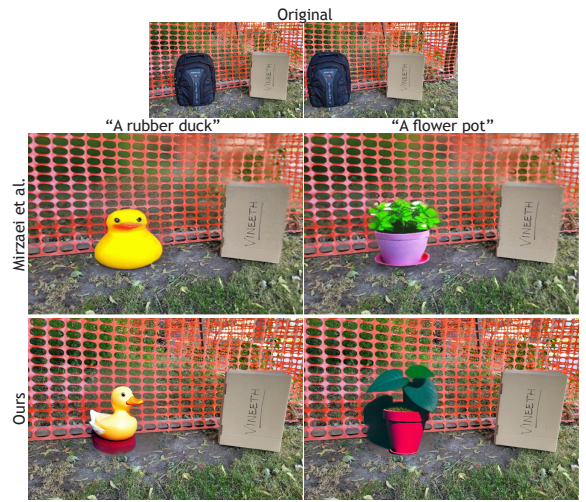


Figure 7. Qualitative comparison with Reference-guided inpainting by [28] for object replacement.

## 5.2. Quantitative Results

3D scene editing is a highly subjective task. Thus, we mainly show various types of qualitative results and comparisons, and recommend readers to refer to the [project page](#) for more results. However, we follow Instruct-NeRF2NeRF and report 2 auxiliary metrics: CLIP Text-Image Direction Similarity and CLIP direction consistency, as shown in Table 1. We compare our method quantitatively with Instruct-NeRF2NeRF and Blended-NeRF for the task of object-replacement on two datasets GARDEN and FERN for various prompts.

Table 1 shows that our method achieves the highest score for CLIP Text-Image Direction Similarity. Interestingly, Blended-NeRF directly optimizes for similarities between CLIP embeddings of the image with the generated object and target text prompts, yet it still achieves a lower score than our method. For Direction Consistency Score,

which measures temporal consistency loss, we observe that Instruct-NeRF2NeRF scores higher than our method on edit prompts where it completely fails (see Figure 5). For example, for the edit "pineapple" in the GARDEN dataset, Instruct-NeRF2NeRF not only fails to create the details of the pineapple but also removes high-frequency details in the background, resulting in a blurry background. We hypothesize that this boosts the consistency score even when the edit is unsuccessful. Therefore, we refer readers to the comparisons in the project video for more details.

## 5.3. Beyond object replacements

**Removing objects** To modify the scene with new contents, ReplaceAnything3D performs objects removal and background inpainting before adding new foreground objects to the scene. Although object removal is not the focus of our work, here we show qualitative comparison with other NeRF-inpainting methods, in particular SPin-NeRF

Table 1. We compute CLIP-based metrics for various datasets: (Top) GARDEN, (Middle) FACE, (Bottom) FERN.

Prompts	CLIP Text-Image Direction Similarity $\uparrow$		CLIP Direction Consistency $\uparrow$	
	Ours	Instruct-NeRF2NeRF	Ours	Instruct-NeRF2NeRF
Pineapple	<b>0.2041</b>	0.0661	0.9590	<b>0.9660</b>
Chess	<b>0.1200</b>	0.0061	0.9457	<b>0.9705</b>
	Ours	BlendedNerf	Ours	BlendedNerf
Mushroom	<b>0.0928</b>	0.0535	<b>0.9781</b>	0.9748
Strawberry	<b>0.3165</b>	0.2224	<b>0.9808</b>	0.9698

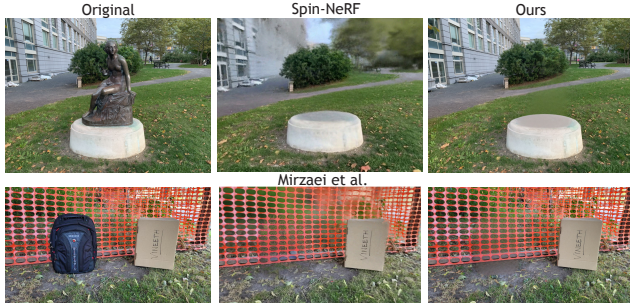


Figure 8. Qualitative comparison for object removal and background inpainting task. Although object removal is not the main focus of ReplaceAnything3D, our method can achieve competitive results with state-of-the-art methods.

[29] and work by Mirzaei et al. [27], to show the effectiveness of our Erase stage. Note that both of these methods only work with forward-facing scenes as shown in Figure 8. Meanwhile, other scene editing technique that works with 360° scenes such as Instruct-NeRF2NeRF is not capable of object removal, as shown in Figure 5.

**Adding objects** In addition to removing and replacing objects in the scene, our method can add new objects based on users’ input masks. Figure 9 demonstrates that completely new objects with realistic lighting and shadows can be generated and composited to the current 3D scene. Notably, as shown in Figure 9-bottom row, our method can add more than one object to the same scene while maintaining realistic scene appearance and multi-view consistency.

#### 5.4. Scene editing with personalized contents

In addition to text prompts, RAM3D enables users to replace or add their own assets to 3D scenes. This is achieved by first fine-tuning a pre-trained inpainting diffusion model with multiple images of a target object using Dreambooth [41]. The resulting fine-tuned model is then integrated into RAM3D to enable object replacement in 3D scenes. As shown in Figure 10, after fine-tuning, RAM3D can effectively replace or add the target object to new 3D scenes.

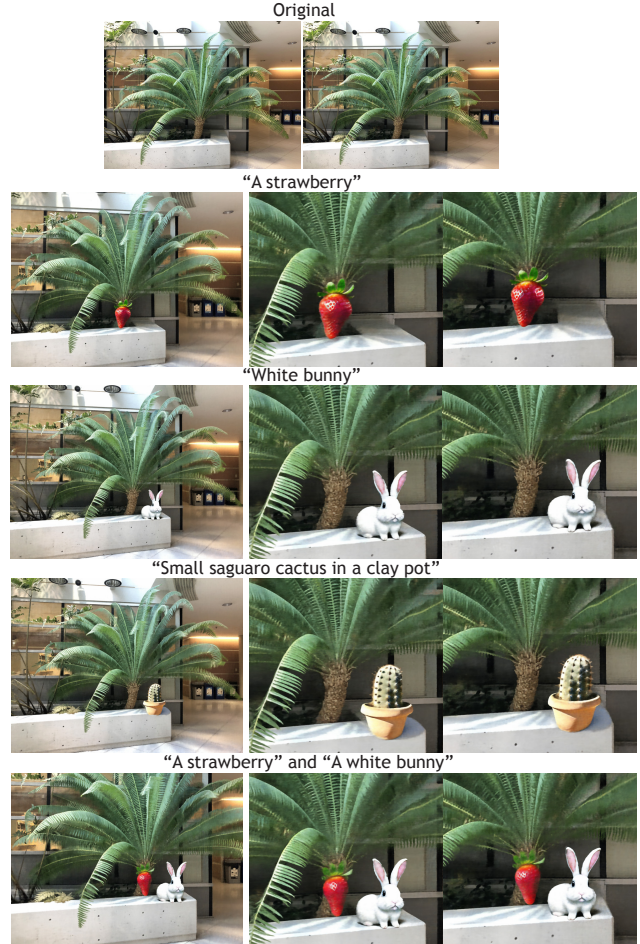


Figure 9. Given user-defined masks, ReplaceAnything3D can add completely new objects that blend in with the rest of the scene. Due to its compositional structure, RAM3D can add multiple objects to 3D scenes while maintaining realistic appearance, lighting, and multi-view consistency (bottom row).

#### 5.5. Ablation studies

We conduct a series of ablation studies to demonstrate the effectiveness of our method and training strategy. In Fig-



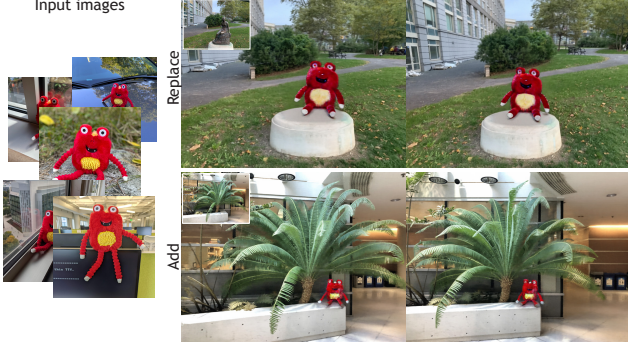


Figure 10. Users can personalize a 3D scene by replacing or adding their own assets using a fine-tuned RAM3D. We achieve this by first fine-tuning an inpainting diffusion model with five images of the target object (left), and then combining it with RAM3D to perform object replacement and addition with custom content.

Figure 11, we show the benefits of our compositional foreground/background structure and background augmentation training strategy. Specifically, we train a version of RAM3D using a monolithic NeRF to model both the background and the new object (combining  $\theta_{bg}$  and  $\theta_{fg}$ ). In other words, this model is trained to edit the scene in one single stage, instead of separate Erase and Replace stages. We observe lower quality background reconstruction in this case, as evident from the blurry hedge behind the corgi’s head in Figure 11-a.

We also demonstrate the advantage of using random background augmentation in separating the foreground object from the background (see Section 4.3). Without this augmentation, the model is unable to accurately separate the foreground and background alpha maps, resulting in a blurry background and floaters that are particularly noticeable when viewed on video (Figure 11-b). In contrast, our full composited model trained with background augmentation successfully separates the foreground and background, producing sharp results for the entire scene (Figure 11-c).

In Fig. 12, we show the importance of the Halo region supervision during the Erase training stage. Without it, our model lacks important nearby spatial information, and thus cannot successfully generate the background scene.

## 5.6. Discussion

Because of our Erase-and-Replace approach for scene editing, our method might remove important structural information from the original objects. Therefore, our method is not suitable for editing tasks that only modify objects’ properties such as their appearance or geometry (for example, turning a bronze statue into a gold one). Furthermore, as ReplaceAnything3D is based on text-to-image model distillation techniques, our method suffers from similar artifacts to these methods, such as the Janus multi-face problem.

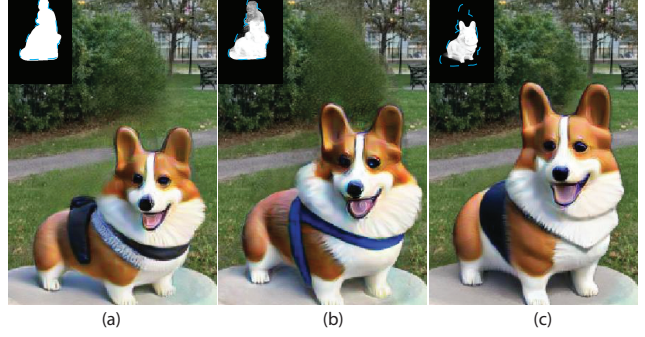


Figure 11. Ablation results for 3 RAM3D variants, on the statue scene for prompt “A corgi”. RGB samples are shown with accumulated NeRF density (alpha map) in the top-left corner. The bubble rendering region is shown as a dotted blue line. a) A monolithic scene representation which contains both the foreground and background. b) A compositional scene model but without random background augmentation. c) Our full model.

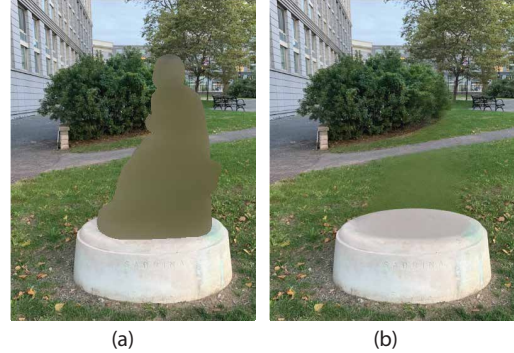


Figure 12. Ablation results for 2 RAM3D variants trained on the Statue scene for the Erase stage. a) Training without any supervision on the halo region surrounding the inpainting mask. The training objective is ambiguous and the Bubble-NeRF model collapses to a hazy cloud. b) Adding halo losses ( $\mathcal{L}_{recon}$  and  $\mathcal{L}_{vgg}$ ) for the halo region surrounding the Bubble-NeRF guides the distillation of  $\theta_{bg}$  towards the true background, as observed on rays which pass nearby to the occluding object. RAM3D can now inpaint the background scene accurately.

In this work, we adopt implicit 3D scene representations such as NeRF or Instant-NGP. For future work, our method can also be extended to other representations such as 3D Gaussian splats [15], similar to DreamGaussian [46]. Interesting future directions to explore include disentangling geometry and appearance to enable more fine-grained control for scene editing, addressing multi-face problems by adopting prompt-debiasing methods [12] or models that are pre-trained on multiview datasets [35, 43], and developing amortized models to speed up the object replacement process, similar to Lorraine et al. [23].



## 6. Conclusion

In this work, we present ReplaceAnything3D, a text-guided 3D scene editing method that enables the replacement of specific objects within a scene. Unlike other methods that modify existing object properties such as geometry or appearance, our Erase-and-Replace approach can effectively replace objects with significantly different contents. Additionally, our method can remove or add new objects while maintaining realistic appearance and multi-view consistency. We demonstrate the effectiveness of ReplaceAnything3D in various realistic 3D scenes, including forward-facing and 360° scenes. Our approach enables seamless object replacement, making it a powerful tool for future applications in VR/MR, gaming, and film production.

## References

- [1] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended Diffusion for Text-driven Editing of Natural Images. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18187–18197, New Orleans, LA, USA, 2022. IEEE. 2
- [2] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended Latent Diffusion, 2023. arXiv:2206.02779 [cs]. 2, 5
- [3] 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 (ICCV)*, pages 5855–5864, 2021. 2
- [4] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022. 5
- [5] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. *ICCV*, 2023. 2
- [6] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. InstructPix2Pix: Learning to Follow Image Editing Instructions, 2023. arXiv:2211.09800 [cs]. 2
- [7] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 3
- [8] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. *arXiv preprint arXiv:2212.08051*, 2022. 3
- [9] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. *arXiv preprint arXiv:2103.10380*, 2021. 2
- [10] Ayaan Haque, Matthew Tancik, Alexei A. Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-NeRF2NeRF: Editing 3D Scenes with Instructions, 2023. arXiv:2303.12789 [cs]. 2, 4, 5
- [11] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-Prompt Image Editing with Cross Attention Control, 2022. arXiv:2208.01626 [cs]. 2
- [12] Susung Hong, Donghoon Ahn, and Seungryong Kim. De-biasing Scores and Prompts of 2D Diffusion for View-consistent Text-to-3D Generation, 2023. arXiv:2303.15413 [cs]. 8
- [13] Clément Jambon, Bernhard Kerbl, Georgios Kopanas, Stavros Diolatzis, Thomas Leimkühler, and George Drettakis. Nerfshop: Interactive editing of neural radiance fields”. *Proceedings of the ACM on Computer Graphics and Interactive Techniques*, 6(1), 2023. 2
- [14] Animesh Karnewar, Tobias Ritschel, Oliver Wang, and Niloy J. Mitra. Relu fields: The little non-linearity that could. *Transactions on Graphics (Proceedings of SIGGRAPH)*, 41(4):13:1–13:8, 2022. 2
- [15] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. 8, 13
- [16] Diederick P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 13
- [17] Quewei Li, Feichao Li, Jie Guo, and Yanwen Guo. Uhdnerf: Ultra-high-definition neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 23097–23108, 2023. 2
- [18] Tianle Li, Max Ku, Cong Wei, and Wenhu Chen. Dreamedit: Subject-driven image editing. *Transactions on Machine Learning Research*, 2023. 2, 11, 13
- [19] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [20] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [21] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object, 2023. 3
- [22] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. Editing conditional radiance fields. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 2
- [23] Jonathan Lorraine, Kevin Xie, Xiaohui Zeng, Chen-Hsuan Lin, Towaki Takikawa, Nicholas Sharp, Tsung-Yi Lin, Ming-Yu Liu, Sanja Fidler, and James Lucas. Att3d: Amortized text-to-3d object synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 17946–17956, 2023. 8
- [24] Luca Medeiros. Language segment anything. GitHub repository, 2021. 2, 4, 11

- [25] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-NeRF for Shape-Guided Generation of 3D Shapes and Textures, 2022. [arXiv:2211.07600 \[cs\]](#). [3](#)
- [26] 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. In *ECCV*, 2020. [2](#), [3](#), [5](#)
- [27] 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 *ICCV*, 2023. [2](#), [7](#)
- [28] 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, 2023. [arXiv:2304.09677 \[cs\]](#). [5](#), [6](#)
- [29] 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, 2023. [arXiv:2211.12254 \[cs\]](#). [2](#), [5](#), [7](#)
- [30] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text Inversion for Editing Real Images using Guided Diffusion Models, 2022. [arXiv:2211.09794 \[cs\]](#). [2](#)
- [31] Chong Mou, Xintao Wang, Jiechong Song, Ying Shan, and Jian Zhang. DragonDiffusion: Enabling Drag-style Manipulation on Diffusion Models, 2023. [arXiv:2307.02421 \[cs\]](#). [2](#)
- [32] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, 2022. [2](#), [3](#), [11](#)
- [33] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. [2](#)
- [34] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. DreamFusion: Text-to-3D using 2D Diffusion, 2022. [arXiv:2209.14988 \[cs, stat\]](#). [3](#), [13](#)
- [35] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, and Bernard Ghanem. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843*, 2023. [8](#)
- [36] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, and Bernard Ghanem. Magic123: One Image to High-Quality 3D Object Generation Using Both 2D and 3D Diffusion Priors, 2023. [arXiv:2306.17843 \[cs\]](#). [3](#)
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Aspell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021. [5](#)
- [38] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *ArXiv*, abs/2204.06125, 2022. [2](#)
- [39] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–12188, 2021. [4](#), [13](#)
- [40] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. [2](#)
- [41] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation, 2023. [arXiv:2208.12242 \[cs\]](#). [7](#)
- [42] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*, 2022. [2](#)
- [43] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. MVDream: Multi-view Diffusion for 3D Generation, 2023. [arXiv:2308.16512 \[cs\]](#). [8](#)
- [44] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. [13](#)
- [45] Hyeonseop Song, Seokhun Choi, Hoseok Do, Chul Lee, and Taehyeong Kim. Blending-NeRF: Text-Driven Localized Editing in Neural Radiance Fields, 2023. [arXiv:2308.11974 \[cs\]](#). [2](#)
- [46] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023. [8](#)
- [47] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-It-3D: High-Fidelity 3D Creation from A Single Image with Diffusion Prior, 2023. [arXiv:2303.14184 \[cs\]](#). [4](#), [13](#)
- [48] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3835–3844, 2022. [2](#)
- [49] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12619–12629, 2023. [3](#)

- [50] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. 11
- [51] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation, 2023. arXiv:2305.16213 [cs]. 3
- [52] Silvan Weder, Guillermo Garcia-Hernando, Áron Monszpart, Marc Pollefeys, Gabriel Brostow, Michael Firman, and Sara Vicente. Removing objects from neural radiance fields. In *CVPR*, 2023. 2
- [53] Jiawei Yang, Marco Pavone, and Yue Wang. 2
- [54] Lin Yen-Chen. Nerf-pytorch. <https://github.com/yenchelin/nerf-pytorch/>, 2020. 11, 13
- [55] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2
- [56] Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yuewen Ma, Rongfei Jia, and Lin Gao. Nerf-editing: geometry editing of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18353–18364, 2022. 2
- [57] Junzhe Zhu and Peiye Zhuang. HiFA: High-fidelity Text-to-3D with Advanced Diffusion Guidance, 2023. arXiv:2305.18766 [cs]. 2, 3, 4, 13
- [58] Jingyu Zhuang, Chen Wang, Lingjie Liu, Liang Lin, and Guanbin Li. DreamEditor: Text-Driven 3D Scene Editing with Neural Fields, 2023. arXiv:2306.13455 [cs]. 2, 4

## A. Additional qualitative comparisons

In Figure 13, we compare our approach with a naive 2D baseline where each image is processed individually. For each image in the training set (first row), we mask out the foreground object (*statue*) and replace it with a new object (*corgi*) using a pre-trained text-to-image inpainting model (Figure 13-second row). We then train a NeRF scene with these modified images. As shown in Figure 13-third row, this results in a corrupted, inconsistent foreground object since each view is very different from each other, in contrast to our multi-view consistent result.

In Figure 14, we demonstrate competitive performance with DreamEditor [18]. It is important to note that DreamEditor has limitations in terms of handling unbounded scenes due to its reliance on object-centric NeUS [50]. Additionally, since DreamEditor relies on mesh representations, it is not clear how this method will perform on editing operations such as object removal, or operations that require significant changes in mesh topologies.

## B. Implementation details

### B.1. NeRF architecture

We use an Instant-NGP [32] based implicit function for the RAM3D NeRF architecture, which includes a memory- and speed-efficient Multiresolution Hash Encoding layer, together with a 3-layer MLP, hidden dimension 64, which maps ray-sample position to RGB and density. We do not use view-direction as a feature. NeRF rendering code is adapted from the nerf-pytorch repo [54].

### B.2. Monolithic vs Erase+Replace RAM3D

We use a 2-stage Erase-and-Replace training schedule for the STATUE, RED-NET and GARDEN scenes. For the FERN scene, we use user-drawn object masks which cover a region of empty space in the scene, therefore object removal is redundant. In this case, we perform object addition by providing the input scene-images as background compositing images to RAM3D.

### B.3. Input Masks

We obtain inpainting masks for object removal by passing dataset images to an off-the-shelf text-to-mask model [24], which we prompt with 1-word descriptions of the foreground objects to remove. The prompts used are: STATUE scene: "statue", GARDEN scene: "Centrepiece", RED-NET scene: "Bag". We dilate the predicted masks to make sure they fully cover the object.

For the Erase experiments, we compute nearby pixels to the exterior of the inpainting mask, and use them as the Halo region (Fig 3). We apply reconstruction supervision on the Halo region as detailed in B.5. For the object-addition





Figure 13. Qualitative comparisons between our method RAM3D (last row) with a naive 2D baseline method, which produces view-inconsistent results (third row). This is because each input image is processed independently and thus vary widely from each other (second row).

experiments in the FERN scene, we create user-annotated masks in a consistent position across the dataset images, covering an unoccupied area of the scene.

#### B.4. Cropping the denoiser inputs

The LDM denoising U-net takes input images of size  $512 \times 512$ . In contrast, RAM3D model outputs are of equal resolution to the input scene images, which can be non-square. To ensure size compatibility, we need to crop and resize the RAM3D outputs to  $512 \times 512$  before passing them to the denoiser (Fig 2). For the STATUE and GARDEN scenes, we resize all images to height 512 and take a centre-crop of  $512 \times 512$ , which always contains the entire object mask re-

gion. For the RED-NET scene, the object mask is positioned on the left side of the images; we therefore select the left-most 512 pixels for cropping.

For the FERN scene, input images are annotated with small user-provided masks. We find that the previous approach provides too small of a mask region to the LDM’s denoiser. In this case, we train RAM3D using the original dataset downsampled by a factor of 2 to a resolution of  $2016 \times 1512$ , and select a rectangular crop around the object mask. We compute the tightest rectangular crop which covers the mask region, and then double the crop-region height and width whilst keeping its centre intact. Finally, we increase the crop region height and width to the max of the

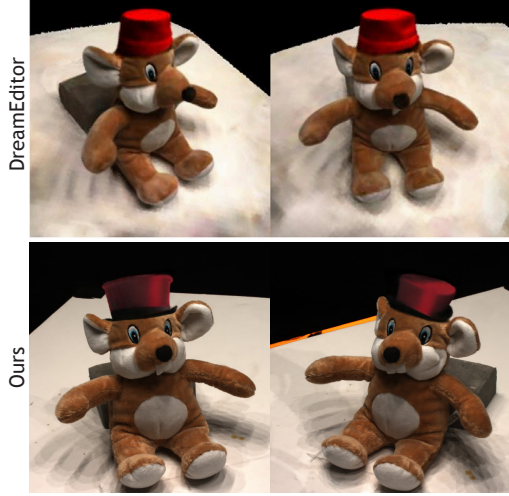


Figure 14. Qualitative comparison with DreamEditor [18] for object addition. Figure adapted from the original DreamEditor paper.

height and width, obtaining a square crop containing the inpainting mask region. We apply this crop to the output of RAM3D and then interpolate to  $512 \times 512$  before proceeding as before.

### B.5. Loss functions

During the Erase training stage, we find it necessary to backpropagate reconstruction loss gradients through pixels close to the inpainting mask (Fig 12), to successfully reconstruct the background scene. We therefore additionally render pixels inside the Halo region (Section B.3, Fig 3), and compute reconstruction loss  $\mathcal{L}_{recon}$  and perceptual loss  $\mathcal{L}_{vgg}$  on these pixels, together with the corresponding region on the input images. Note that the masked image content does not fall inside the Halo region in the input images - therefore  $\mathcal{L}_{recon}$  and  $\mathcal{L}_{vgg}$  only provide supervision on the scene backgrounds. For the reconstruction loss, we use mean-squared error computed between the input image and RAM3D’s RGB output. For perceptual loss, we use mean-squared error between the features computed at layer 8 of a pre-trained and frozen VGG-16 network [44]. In both cases, the loss is calculated on the exterior of the inpainting mask and backpropagated through the Halo region. During the Replace training phase, following Zhu and Zhuang [57], we apply  $\mathcal{L}_{BGT+}$  loss between our rendered output  $\mathbf{x}$ , and the LDM denoised output  $\hat{\mathbf{x}}$ , obtaining gradients to update our NeRF-scene weights towards the LDM image prior (see HiFA Loss in Fig 2, eqn 11 [57]). No other loss functions are applied during this phase, thus loss gradients are only backpropagated to the pixels on the interior of the inpainting masks. For memory and speed efficiency, RAM3D only renders pixels which lie inside the inpainting mask at this stage (Fig 4), and otherwise samples RGB values directly

from the corresponding input image.

Finally, following Tang et al. [47], we apply depth regularisation using the negative Pearson correlation coefficient between our NeRF-rendered depth map, and a monocular depth estimate computed on the LDM-denoised RGB output. The depth estimate is obtained using an off-the-shelf model [39]. This loss is backpropagated through all rendered pixels; i.e the union of the inpainting mask and Halo region shown in Fig 3. We do not apply this regularisation during the Replace stage. In summary, the total loss function for the Replace stage is:

$$\mathcal{L}_{total} = \mathcal{L}_{BGT+} + \lambda_{depth} \mathcal{L}_{depth} + \lambda_{recon} \mathcal{L}_{recon} + \lambda_{vgg} \mathcal{L}_{vgg} \quad (7)$$

with loss weights as follows:  $\lambda_{recon} = 3, \lambda_{vgg} = 0.03, \lambda_{depth} = 3$ .

We use the Adam optimiser [16] with a learning rate of  $1e-3$ , which is scaled up by 10 for the Instant-NGP hash encoding parameters.

### B.6. Other Training details

Following [34, 57], we find that classifier-free guidance (CFG) is critical to obtaining effective gradients for distillation sampling from the LDM denoiser. We use a CFG scale of 30 during the Replace stage, and 7.5 during the Erase stage. We also adopt the HiFA noise-level schedule, with  $t_{min} = 0.2, t_{max} = 0.98$ , and use stochasticity hyperparameter  $\eta = 0$ . In the definition of  $\mathcal{L}_{BGT+}$  loss (see eqn 11 in [57]), we follow HiFA and choose a  $\lambda_{rgb}$  value of 0.1. We render the RAM3D radiance function using a coarse-to-fine sampling strategy, with 128 coarse and 128 fine raysamples. During the Replace training stage, we swap the composited background image with a randomly chosen plain RGB image at every 3rd training step. As shown in Fig 11, this step is critical to achieving a clean separation of foreground and background.

We train RAM3D for 20,000 training steps, during both Erase and Replace training stages, which takes approximately 12 hours on a single 32GB V100 GPU. The output of Replace stage training is a set of multiview images which match the input scene images on the visible region, and contain inpainted content on the interior of the masked region which is consistent across views. To obtain novel views, we train standard novel view synthesis methods using RAM3D edited images and the original scene cameras poses as training datasets. We use nerf-pytorch [54] for the LLFF scenes (STATUE, FERN, RED-NET SCENES), and Gaussian Splatting [15] for the GARDEN scene.