

Language-driven Object Fusion into Neural Radiance Fields with Pose-Conditioned Dataset Updates

Ka Chun Shum¹ Jaeyeon Kim¹ Binh-Son Hua² Duc Thanh Nguyen³ Sai-Kit Yeung¹

¹Hong Kong University of Science and Technology ²Trinity College Dublin ³Deakin University

Abstract

Neural radiance field is an emerging rendering method that generates high-quality multi-view consistent images from a neural scene representation and volume rendering. Although neural radiance field-based techniques are robust for scene reconstruction, their ability to add or remove objects remains limited. This paper proposes a new language-driven approach for object manipulation with neural radiance fields through dataset updates. Specifically, to insert a new foreground object represented by a set of multi-view images into a background radiance field, we use a text-to-image diffusion model to learn and generate combined images that fuse the object of interest into the given background across views. These combined images are then used for refining the background radiance field so that we can render view-consistent images containing both the object and the background. To ensure view consistency, we propose a dataset updates strategy that prioritizes radiance field training with camera views close to the already-trained views prior to propagating the training to remaining views. We show that under the same dataset updates strategy, we can easily adapt our method for object insertion using data from text-to-3D models as well as object removal. Experimental results show that our method generates photorealistic images of the edited scenes, and outperforms state-of-the-art methods in 3D reconstruction and neural radiance field blending.

1. Introduction

Editing of 3D scenes by insertion or removal of objects has been a fundamental task in computer graphics and computer vision, which can be achieved using traditional 3D scene authoring tools [4, 8]. For example, to insert an object into a 3D scene, traditional approach requires user to manually select the object and position it into the scene. This manual pipeline has been used in a wide range of applications, such as furniture arrangement in interior design [14]

and asset creation in game development [9, 12].

Recent advances in deep learning for image synthesis have opened new directions to scene editing. Neural radiance field (NeRF) [39] is a pioneering method that can learn to render view-consistent photorealistic images using neural networks. Generative models such as generative adversarial networks (GANs) [11] and diffusion models [16] learn to output photorealistic images from unconstrained image collections. Moreover, recent text-guided diffusion models [49, 50, 52] show great promises in generating high-quality realistic and diverse images from a single text prompt. Such capability of language-driven image synthesis inspires us to revisit 3D scene editing using natural languages.

In this paper, we propose a new language-driven method for editing neural radiance fields and manipulating their objects. Particularly, we focus on the task of inserting a foreground object into a background radiance field. Our approach to this problem is to utilize a generative model to synthesize new images containing both the given object and background. The synthesized images can be subsequently used to refine the background radiance field to learn new object geometry and appearance. This approach to refining the radiance field is also known as *dataset updates* [13]. However, a notable challenge associated to dataset updates is that the refining process might degrade the rendering quality due to inconsistent views synthesized by the generative model, and thus affecting the geometry, appearance, and convergence of the background radiance field. To address this issue, we propose a new strategy for dataset updates that controls the refinement process so that the foreground object is gradually introduced into the training, beginning at a randomly selected view and then prioritizing views close to the already-trained camera views before propagating the refinement to views further away. We observe that this strategy improves the learning of the background radiance field in integrating the foreground object, and significantly reduces rendering artifacts while maintaining view-consistent rendering. We show that this *pose-conditioned* strategy is robust to object insertion for both real objects and virtual

objects from recent text-to-3D models [30, 46, 63]. More importantly, we show that under this strategy, we can also refine the background radiance field for object removal.

In summary, our contributions are as follows.

- A new framework based on neural radiance fields and text-to-image diffusion for object insertion and removal in 3D scenes;
- A pose-conditioned dataset updates strategy that stabilizes the refinement of the background radiance field in the presence of a new foreground object, resulting in a view-consistent object fusion;
- A set of extensive experiments that validate the robustness of our method and demonstrate its state-of-the-art performance as well as an ablation study to analyze different factors of our proposed method.

2. Related Works

Object manipulation has been a long-standing research problem in computer graphics and computer vision. Here, we limit our discussion to neural network-based methods with a focus on generative models and neural radiance fields.

Image Manipulation. In the era of deep learning, image manipulation can be achieved via conditional generation by popular generative models such as GANs [11] and diffusion models [16]. Early methods perform conditional data generation by combining supervised losses with adversarial losses to learn class-conditional adversarial networks [40] and image-to-image translations [19, 62, 75]. Subsequently, StyleGAN [21] and its variances [20, 22] proposed to perform image editing by traversing their semantic latent representation, such as interpolation in the latent space, but such conditioning remains implicit and difficult to control.

The introduction of diffusion models [16] has shown great potential in image synthesis with high-quality data samples constructed from its sophisticated forward and denoising diffusion steps. Recent developments of vision-language models such as CLIP [48] have led to the popularity of using text prompts as an intuitive condition to generate and edit images. Several text-guided diffusion models [49, 50, 52] offer excellent image quality from training with extraordinary data and computational resources. Applications can be built upon these methods by fine-tuning these text-guided models for downstream tasks [3, 23, 35, 51, 66, 73]. For example, RePaint [35] and SmartBrush [66] inpaint images with masks by denoising. ControlNet [73] fine-tunes a twin diffusion model that accepts custom input images. Instruct-Pix2Pix [3] and Imagic [23] retrain on an edit-text-to-image dataset to allow text-guided editing. Our method is part of this line of work where we use a text-to-image model [51] to guide the fusion of new objects into a background.

Neural 3D modeling. Early image-based 3D modeling methods using convolutional neural networks (CNNs) [27] simply learn by stacking multiple images as input or output [10, 17, 32, 57, 74], where CNNs struggle to deal with complex shapes, textures, and lighting implicitly captured in these images. Follow-up works integrate differentiable rendering and represent 3D structures as neural surfaces [38] or shapes [7]. 3D-aware GANs integrates volume rendering as part of their generators to synthesize novel views from a single image [5, 44, 55, 56]. Neural radiance fields (NeRF) [39] have a similar neural rendering approach, but are optimized on a ray rendering loss on multi-view input images. A family of subsequent works addresses the limitation of the original NeRFs in the perspective of visual artifacts [1], data complexity [36], camera poses [70], or computational efficiency [42, 72].

NeRF editing. Editing of a 3D scene using a NeRF representation has recently received considerable attention. A basic approach is to directly edit a pre-trained NeRF by several methods including parameter tuning [34], layer feature fusion [58], or deformable rays [71]. Another manner is to directly manipulate the multi-view images and the training process of NeRFs. For example, some NeRF stylization methods freeze the geometry branch and optimize the color branch in a NeRF to stylize multi-view images [18, 43, 45, 59]. Several methods attempt to decompose items in existing NeRFs, which in turn hold the colors and separate voxels considering multi-view masks [28] or semantics [25, 26]. NeRF inpainting fills simple unseen background geometry and colors with help of depth priors [41] or filtering inpainted multi-views [64].

Creating simultaneously complex geometry and vivid color in a pre-trained NeRF is much more challenging due to higher-level requirements of consistency. DiscoScene [68] fuses results from two random background and object NeRF generators, thus is unable to condition any customized content. BlendNeRF [25] automatically aligns poses and blends two NeRFs for human or animal face images. FocalDreamer [29] and DreamEditor [76] rely on fine meshes to function, disregarding the advantage of multi-view representation of NeRFs. Our method is perhaps most similar to Instruct-NeRF2NeRF [13] which shares a related data updating schema. However, compared to our method, Instruct-NeRF2NeRF [13] always generates geometry aligned with the original geometry (e.g., transfer a sneaker to a sneaker-shape apple), thus failing in most cases of object insertion. Also, it requires heavy retraining on a large-scale self-constructed dataset.

Recently, there is an emerging trend of using natural languages to generate 3D models from priors guided by text-to-image diffusion. DreamFusion [46] introduced score distillation sampling (SDS) loss that progressively consolidates the view information from a diffusion model into

NeRFs. This loss has been applied with additional treatments on image resolution [30], conditional images [37], photo-realism [63], or scene geometry [67]. We adopt the SDS loss but develop a novel training schedule to address challenges in multi-view object and background fusion for neural radiance fields editing.

3. Method

3.1. Overview

We first describe our method for the task of object insertion into a background NeRF. We will then show that our method generalizes to object removal under the same principle. An overview of our method is illustrated in Fig. 1.

Our method accomplishes object insertion in two steps. In the first step, given a target object and a radiance field, we synthesize training views (images) for the radiance field with the target object (i.e., the target object is included in the training views). In the second step, the radiance field is refined by training with the new training views to learn the geometry and appearance of the new object. The training views in the first step can be created using generative image synthesis. We expect the refinement of the radiance field to learn 3D representations consistent with input views from the first step. However, there are two challenges. First, the background is required to be preserved in the synthesized images so that the subsequent updating step results in a radiance field with the original background. Second, updating of the radiance field with the synthesized images may result in artifacts due to inconsistent input views generated from image synthesis. Our solution to these challenges to achieve high-quality object insertion is as follows.

First, we leverage a state-of-the-art text-to-image diffusion model for our image synthesis, and opt to customize the model for our background preservation purpose. Particularly, we fine-tune a pre-trained text-to-image diffusion model with the target object and background (from the radiance field) so that the model can generate new object images with background preserving. We formulate this task as an image inpainting problem. Specifically, we use a binary mask to indicate the region in a background image that the object should be inserted. Therefore, performing image inpainting with content generation in the mask region allows object insertion. To achieve this inpainting capability, we aim to train a diffusion model to complete the background with a desirable object inside the mask. Here, we design our diffusion model based on the inpainting network of Stable Diffusion [50]. We further build upon the interesting ability of DreamBooth [51] that personalizes the diffusion model to be able to generate images containing our object of interest and background. This personalization is described by text prompts that contain special *identifier* tokens for the object of interest and the background, respectively. Unfortunately,

this approach only works well for single-image object insertion. Performing background for multiple camera views results in inconsistent results, making the training of the radiance field fail to converge to desirable quality.

To achieve accurate and consistent synthesis of both object and background in multiple views, we propose a new strategy to schedule the data used in the training of the radiance field during refinement. Our strategy is inspired by an important observation about the nature of NeRF: a view rendered by a NeRF maintains an extent of color information from nearby already-trained views, the nearer the more noticeable. Therefore, if we pass these nearby renderings to the fine-tuned diffusion model with properly controlled noises, view-consistent results can be generated based on the learned color hints. Reversely, the outputs from the fine-tuned diffusion model can be added into the training data of the NeRF. As a consequence, the prior object knowledge from the diffusion model in 2D is gradually consolidated into the background NeRF in 3D.

Based on such observation, we design a novel data scheduler for our NeRF training as follows. We first optimize the NeRF by following the traditional training procedure on a dataset of multi-view background images. We then progressively fuse the target object into the background NeRF by iteratively updating the training dataset in a pose-ordered manner. We propose to sort the views based on the increasing distance from the camera poses of not-yet-included views to the poses of already-trained views. Therefore, we term this strategy as *pose-conditioned* dataset updates. Starting from an original training dataset, new training views (with the target object inserted in) are generated by the fine-tuned diffusion model. The schedule follows a pose-conditioned rule, for example, views with more similar poses to already-trained views are prioritized for training. Like [13], we also periodically replace some old views in the dataset with the more updated ones, similarly through the fine-tuned diffusion model. The NeRF refinement ends when the images from all camera poses are included in the training.

3.2. Background and Object Fine-tuning

We present our background radiance field by a set of scene images and their camera poses $\{I^b, T\}$. The object of interest is defined by a set of images I^o that captures the object from different viewpoints. Our final goal is to generate a set of view-consistent images \hat{I} that contains the same background from I^b and the object defined by I^o at a location defined by a 3D bounding box B . Note that in our method, I^o can be real-world photographs captured by users, or synthetically produced by off-the-shelf text-to-3D generation methods.

Let D_θ be our diffusion model, which is based on the inpainting model from a pre-trained Stable Diffusion [50]. To personalize the diffusion model on the background images

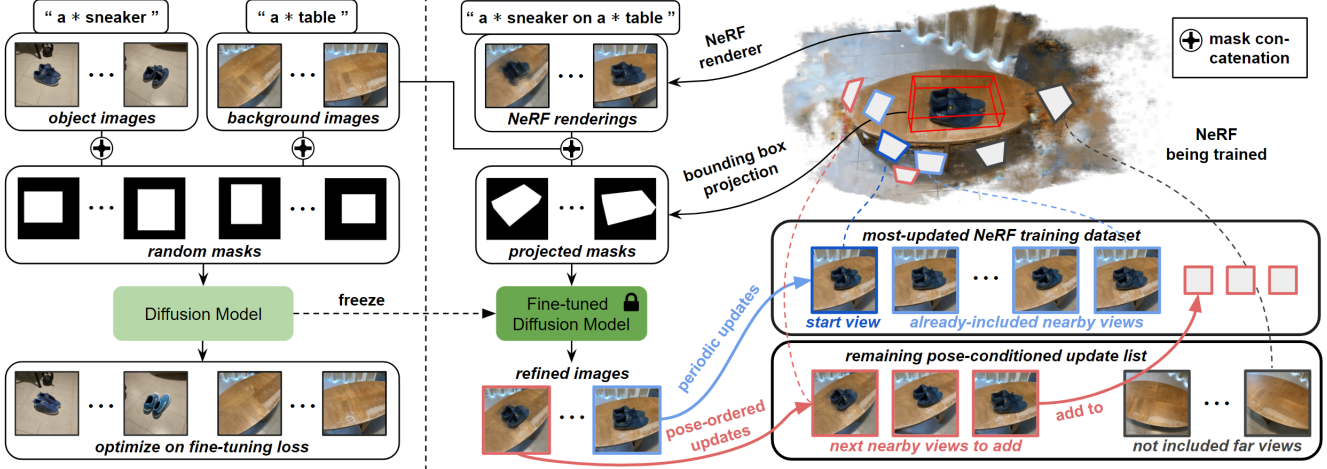


Figure 1. Overview of our pipeline. We first adopt a diffusion model for view synthesis with text identifiers based on DreamBooth [51] but in an inpainting manner (left). Then, we freeze the diffusion model during our pose-conditioned NeRF optimization (right). Information from the diffusion model and the NeRF iterates mutually. Views generated by the diffusion model are added to our on-going dataset to strengthen the NeRF 3D. In return, the NeRF renders color hints for diffusion model to refine new views.

I^b and object images I^o , we learn to predict the noise:

$$\hat{\mathcal{N}}(t) = D_\theta(I^k, M, \mathcal{N}(t), P^k), \quad (1)$$

where θ represents the parameters of the diffusion model D and $k \in \{b, o\}$ denotes background and object, respectively. $M \sim \mathcal{U}$ is a random square mask sampled from a uniform distribution \mathcal{U} , that is, the coordinates of the four mask corners are uniformly sampled in a value range. $\mathcal{N}(t)$ refers to the process of sampling a noise map from a Gaussian distribution \mathcal{N} given a random noise strength $t \in [0, 1]$. P^b and P^o are the background prompt and object prompt containing the identifier word as in DreamBooth [51].

To fine-tune the diffusion model, we use custom text-image pairs that include a user-defined token to identify our background and object, e.g., we use “*” symbol as an adjective to modify a sneaker noun, meaning that an “* sneaker” is our sneaker of interest. Such text prompt is then paired with a few images of the same object (e.g., photos of the same sneaker) for fine-tuning the diffusion model. The fine-tuned text-to-image model can then generate our background and object instance with unseen camera poses.

Our fine-tuning loss $\mathcal{L}_{\text{finetune}}$ follows the standard reconstruction loss on the predicted noise $\hat{\mathcal{N}}(t)$ and the sampled noise $\mathcal{N}(t)$:

$$\mathcal{L}_{\text{finetune}} = \|\hat{\mathcal{N}}(t) - \mathcal{N}(t)\|^2 \quad (2)$$

We fine-tune D_θ on I^b and I^o with the components mentioned in Eqs. 1 and 2 for n_{bg} and n_{obj} times, respectively. Let the final fine-tuned diffusion model be \hat{D} , which we will use for the subsequent NeRF refinement stage in Sec. 3.3.

Note that Eq. 1 is an abstract representation of our fine-tuning process. Some intermediate steps such as noise

scheduling, noisy image construction, image masking, latent encoding-decoding, and input concatenation have been omitted for simplicity. Please refer to our Appendix for detailed steps.

3.3. Pose-conditioned Dataset Updates

Let us describe our dataset updates rule as follows. Let R_θ be our NeRF to refine, which takes as input a background camera pose T and outputs a rendering $R_\theta(T)$ at that viewpoint. Let image at viewpoint T in the current training dataset as $\hat{I}(T)$, our NeRF training loss $\mathcal{L}_{\text{NeRF}}$ minimizes the difference between NeRF rendering $R_\theta(T)$ and the training data $\hat{I}(T)$:

$$\mathcal{L}_{\text{NeRF}} = \|R_\theta(T) - \hat{I}(T)\|^2. \quad (3)$$

At start, we train our NeRF R_θ on background data $\{I^b, T\}$, and thus $\hat{I}(T)$ is equal to I^b at viewpoint T . Then, we progressively introduce the target object into our NeRF training as follows. Suppose that the camera views already used for our NeRF refinement are represented by a set of poses $\mathbf{T}_i = \{T_0, T_1, \dots, T_i\} \subseteq T$ where T is the set of all poses. To select a new view T_{i+1} , we compute the distance from each candidate view to the existing set of poses as follows:

$$T_{i+1} = \arg \min_{T_j \in T \setminus \mathbf{T}_i} \min_{T_k \in \mathbf{T}_i} \|\text{trans}(T_j) - \text{trans}(T_k)\| \quad (4)$$

where $\text{trans}()$ extracts the translation component of a camera pose. Note that we empirically found that measuring only the Euclidean distance between pose translations is enough for reasonable view selection. It is because our

multi-view training data are collected with a smoothly connected, inward-surrounding, and non-repeated camera trajectory. This requirement is fulfilled by most NeRF dataset and common object insertion scenarios.

Next, let us describe how we generate images for our dataset updates. We generate the starting view $\hat{I}(T_0)$ as:

$$\hat{I}(T_0) = \hat{D}(I^b(T_0) \oplus R_\theta(T_0), B(T_0), t = 1, P^{o+b}) \quad (5)$$

where $B(T)$ refers to the target object location defined by an 2D image mask projected by pose T from the 3D bounding box B . We use symbol \oplus to denote the process of mask-based combination of bounding box $B(T)$ over the ground-truth background $I^b(T_0)$ and NeRF rendering $R_\theta(T_0)$. Content of $I^b(T_0)$ and $R_\theta(T_0)$ are concatenated to the unmasked and masked region, respectively. We fix the noise strength to $t = 1$ to ensure \hat{D} gives good initialization from pure noise. Let P^{o+b} be the target prompt that describes the edited scenes. It is the combination of prompts P^o and P^b but with a preposition in-between (e.g., “a * sneaker on a * table”). P^{o+b} activates \hat{D} to produce content that properly combines the background and object. We maintain the simplicity of our text prompt format for the ease of usage. We also found additional description (e.g, a long adjective) not particularly helpful for the fine-tuning discussed in Sec. 3.2 and the diffusion model inferencing in Eq. 5. It is possibly due to diffusion models that mainly learn the object of interest through the given images and prior knowledge of the object category [51].

As training goes on, more views are sequentially included. In addition to pose-conditioned views, we also periodically replace existing views in the dataset similarly through the diffusion model to fix the minor inconsistency [13]. The newly added views or the periodically updated old views, both represented as $\hat{I}(T_i)$, are generated similarly through the diffusion model as in Eq. 5:

$$\hat{I}(T_i) = \hat{D}(I^b(T_i) \oplus R_\theta(T_i), B(T_i), t = \tau, P^{o+b}) \quad (6)$$

where the noise strength $t = \tau$ is controlled at a relatively low value to allow \hat{D} utilizes the color hints learned from the previous nearby views, and still give enough amount of randomness to fix the geometry and lighting defects. In our implementation, we use $\tau = 0.35$.

In our pose-conditioned dataset updates, we include in total n_{near} new nearest views into the ongoing dataset every n_{new} NeRF training steps. For periodic updating on old views, similar to Instruct-NeRF2NeRF [13], we update one random old view every n_{old} NeRF training steps. Our pipeline stops training until the last view is included and trained another n_{new} steps.

3.4. Integrating Text-to-3D Object

Our method can also support text-driven object insertion using data from synthetic 3D models. With a target

object prompt, we use a text-to-3D method such as DreamFusion [46] to generate a NeRF that represents a 3D object, which we can perform rendering to generate multi-view images for the object I^o , which can be used as input to our pipeline. From our experience, having a number of fixed camera trajectories inward-surrounding the object is sufficient to construct a valid I^o for most objects.

3.5. Adapting to Object Removal Task

Under the same framework, it is possible to adapt our method to perform object removal as follows. As object is now irrelevant due to removal, we do not need to perform fine-tuning of the diffusion model on the foreground object. For the background, we inpaint the masked area of all multi-view background images using the background prompt without the identifier (e.g., “a table”). These inpainted images are treated as *pseudo-ground truth* for fine-tuning the diffusion model. Note that these pseudo backgrounds are visually pleasing but remains inconsistent across views. We can then perform the NeRF refinement step by using iteratively dataset updates, but with the difference of only using the background prompt. We found that this adapted pipeline gradually transforms the masked areas to a consistent background across views. We will analyze in detail in our experiment the necessity of the pseudo background and show what background inconsistency our adapted pipeline fixes.

4. Experiment

4.1. Dataset

We propose a dataset of six scenes that includes multi-view images of various backgrounds and objects. The number of multi-view images of each background and object ranges from 40 to 80 and 20 to 40, respectively. We collect the images with an iPhone SE 2 front camera in resolution 3072x3072 at RAW format. We resize them to 512x512 PNG format for the ease of training.

For object removal, we train and evaluate on the commonly used inpainting dataset from Mip-NeRF-360 [2] and IBRNet [61]. We center crop the interest area of the high-resolution wide-angle images to resolution 512x512 to fit our pipeline.

4.2. Baselines

For comparisons, we select SOTA baselines from different branches of work that can perform object insertion. To ensure fairness, we barely make changes to the baselines, only when necessary to adopt them to our task.

Traditional 3D pipeline. From multi-view images, the mature traditional 3D reconstruction pipeline is able to build a corresponding 3D model. The scene editing is then processed on these 3D models directly. To simulate a traditional computer graphics baseline, we use COLMAP [53, 54] to

reconstruct the mesh of the background and object, and then manually crop and place the object mesh into the background mesh. To ensure minimal of changes, we do not perform additional manual rendering of the texture and lighting.

Image inpainting. Filling content to the view-consistent masked areas of the images is a possible way to edit a multi-view scene. However, the SOTA multi-view consistent inpainting methods only apply to inpaint background for object removal task [41, 60, 64]. An inpainting variant of Stable Diffusion [50] fills up the masked area of a single view with the object specified by a text prompt. We still treat it as a baseline since it generates related objects in our background given our text prompts. For comparison, we perform single-view inpainting on all of our masked background images.

Single-image-to-3D NeRF. Synthesizing NeRF-based novel views from a single image partially fulfills our task. The NeRF learns external knowledge from different prior modules such as CNNs [70], vision transformer [31], and diffusion model [6]. As diffusion model receiving increasing attention for its generation ability and to best match our pipeline, we select a SOTA baseline Zero123 [33] that uses a distillation prior [46] from a 2D diffusion model to synthesize novel views for comparison. Similarly, we use the same starting view in Eq. 5 as the single image.

Instruct-NeRF2NeRF. Instruct-NeRF2NeRF [13] shares a similar NeRF training and dataset updating modules as ours. Although it fails to add or remove objects with noticeable geometry changes, it still injects new 3D information partially from the text prompt. We follow its default setting to train it as a baseline, except for adjusting our text prompt to its required format. For example, the text prompt is changed from “a * sneaker on a * table” to “add a sneaker on a table”.

Other methods. Given multi-view object and background images as input, training two corresponding NeRFs and then blending them into one is possible for scene editing. We do not consider naive and manual blending, as it is similar to the traditional computer graphics pipeline, which we already compared. BlendNeRF [25] automatically combines two NeRFs, but as this method only supports limited image domains, it is not applicable to our diverse scene data for comparison. DreamEditor [76] utilizes DreamBooth [51] for fine-grained editing but this method still fails to construct objects with obvious new geometry.

4.3. Implementation

Below configuration works for most of our editing cases. We fine-tune on object images more times than background images as we found the diffusion model needs extra training to understand the complicated object other than the relatively simple background. We set $n_{bg} = 400$ and $n_{obj} = 4000$. For NeRF training, we empirically found $n_{near} = 3$, $n_{new} = 500$, and $n_{old} = 10$ balance well the quality and efficiency. We fix the noise strength as mentioned in Eqs. 5 and 6.

We run all experiments in image resolution 512x512 on one Nvidia RTX 3090 GPU. Diffusion model fine-tuning in total takes around 40 minutes. NeRF optimization time depends on the number of background images. Every diffusion model inference cost 8 seconds and every NeRF backpropagation cost 0.5 second. Following our configuration, our NeRF costs around 2.7 times more training time than the usual NeRF for the same number of NeRF training steps.

4.4. Qualitative Results

Baseline comparison. We qualitatively compare our method with other baselines in Fig. 2. As shown in the results, our method generates plausible object content in the input background consistently across views. Also, our output images well match the target text prompt and fit the 3D bounding box defined by users.

For the baselines, the traditional 3D pipeline suffers from seamed object-background boundary and unrealistic lighting. It can be seen that geometry in the obscure regions such as corners is not accurately reconstructed. The image inpainting baseline produces reasonable individual images but the generated object is inconsistent in different views. The model has difficulty referring to the same object even using a more detailed target text prompt. Single-image-to-3D NeRF fails on challenging input with complex object and background. A larger camera pose shift can result in ground-truth background mismatch or even content collapse. Instruct-NeRF2NeRF is known to be strong at stylizing existing geometry but weak at generating new geometry. Under our setting, it performs random view-consistent editing, but fails to accomplish object addition task.

Text-to-3D edits. We present our results of text-to-3D object insertion. As shown in Fig. 3, although the text-to-3D objects are synthetic, our method can still produce high-quality results with consistent camera views. We hypothesize that with future improvement of the text-to-3D generation models, the texture and lighting on the 3D model can be improved to better fit the background.

Object removal. We show the object removal results of our adapted pipeline in Fig. 4. Our method generates multi-view consistent background inpainting results on regular cases. Note that lacking the pseudo background fine-tuning module can lead to gradual background collapse, which we will discuss in the ablation study in Sec. 4.6.

4.5. Quantitative Results

Although the edits lack ground truth for comparison, we still use following metrics to evaluate our method. We use CLIP Score [15] to measure how well an image correlates to a target prompt in the CLIP space. We average a score on all edited images \hat{I} paired with their target prompts P^{o+b} . Denote $E(\{I, P\})$ as the CLIP embedding of an image or text

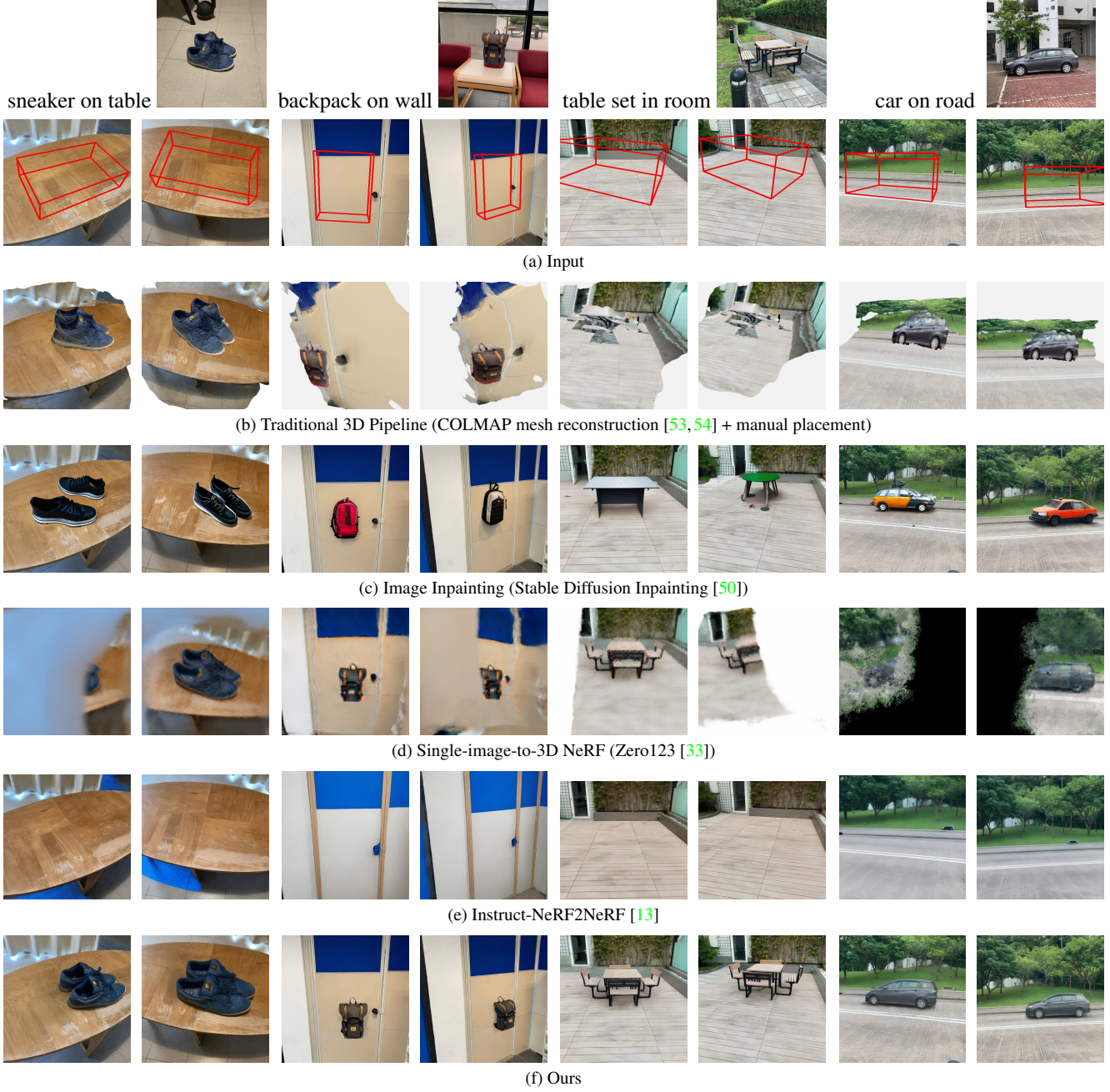


Figure 2. Visualization of the edited scenes. The input of our method includes the multi-view background images and a target 3D bounding box (second row), as well as the multi-view object images and an target text prompt (first row). Note that the baselines may not use all of the input due to the nature of their techniques. More results are in the supplementary material.

prompt, the CLIP Score is formulated as a cosine similarity:

$$\text{CLIPScore} = \cos^+ \left(E(\hat{I}), E(P^{o+b}) \right). \quad (7)$$

where $\cos^+(a, b) = \max(0, \cos(a, b))$. We further use CLIP Directional Consistency [13] to measure the editing quality and consistency across view points in the CLIP space. Specif-

ically, the CLIP space changes from the background I^b to the edited images \hat{I} should match the changes in the prompts P^b and P^{o+b} . Moreover, a consistent 3D edit infers to the consistent changes from every pair of I^b and \hat{I} , especially for adjacent views. We follow Instruct-NeRF2NeRF [13] to formulate the average performance on all adjacent views I_i

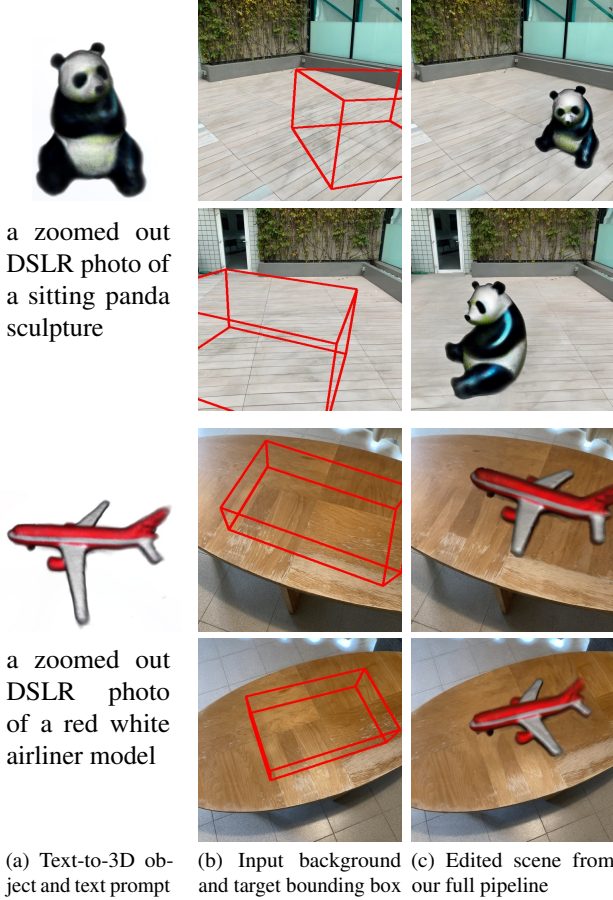


Figure 3. Qualitative results of using text-to-3D objects. After object fine-tuning on the rendered images of the text-to-3D objects (a), the background (b) is filled by the learned objects. Our method generates pose consistent editing results (c) even the text-to-3D objects look synthetic.

and I_{i+1} :

$$\begin{aligned} \text{CLIPDirectionalConsistency} = \\ \cos^+ \left(E(P^{o+b}) - E(P^b), E(\hat{I}_i) - E(I_i^b) \right) \times \\ \cos^+ \left(E(\hat{I}_i) - E(I_i^b), E(\hat{I}_{i+1}) - E(I_{i+1}^b) \right) \end{aligned} \quad (8)$$

The quantitative comparison is shown in Tab. 1 based on the evaluation of six scenes in our dataset. The results show that our method outperforms the baselines especially in terms of the metric that measures the view consistency.

4.6. Ablation Study

We show the effectiveness of our proposed modules. In Fig. 5, two views of an edited scene are shown. Temporal results of these two views are presented to show the gradual failure during the training without specific modules

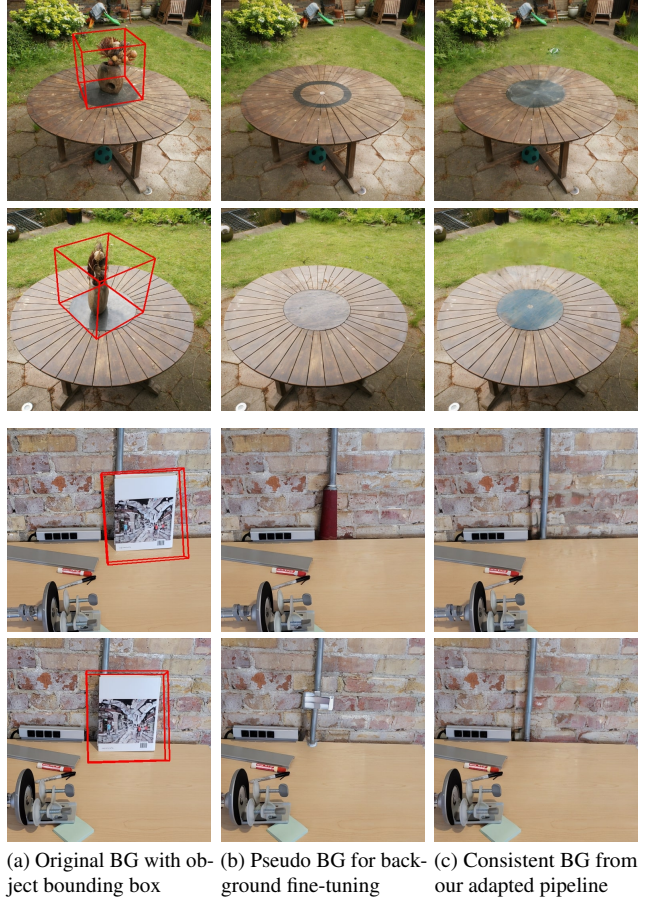


Figure 4. Qualitative results of our adapted pipeline on object removal task. The input images with masks defined by a 3D bounding box (a) are inpainted inconsistently as pseudo background (b) for background fine-tuning. Following our pipeline, our method still converges to a consistent background across views (c).

Method \ CLIP Metric	CLIP Score \uparrow	Directional Consistency \uparrow
Traditional 3D [53, 54]	0.2589	0.0962
Inpainting [50]	0.2665	0.0803
Zero123 [33]	0.2302	0.0399
Instruct-NeRF2NeRF [13]	0.2351	0.0248
Ours	0.2764	0.1845

Table 1. Quantitative results. Higher CLIP metrics scores indicate higher image editing quality or consistency.

in our pipeline. We suggest readers observe the images with the analysis described below.

Object fine-tuning. Fine-tuning on object images keep diffusion model aware of a same object instance during multiple times of model inferences in the dataset updating steps. Without object fine-tuning, object in different views

diverge to different appearance. A not fine-tuned diffusion model only has intent to generate random suitable object but not to target a specific instance.

Background fine-tuning. Except the generated objects, the rest of the masked area filled by the diffusion model should match the ground-truth background. We empirically found that our pipeline runs under a low inference noise strength of the diffusion model (in Eq. 6) causes gradual background collapse. The filled background is darken and clustered with increasing training steps. We hypothesize the reason being diffusion model implicitly keeps some dark noisy pixels as part of its output. After thousands of repeated model inferences in our pipeline, the slightly darken background accumulates to the final collapse.

Fine-tuning on background images helps the diffusion model revise the wrong pixels to the learned background. Note that a high noise strength setting partially solve this issue but causes another consistency problem, which we will discuss in a later subsection in this ablation study.

Pose-conditioned dataset updates. In this experiment, we update and optimize all views in the NeRF training dataset randomly and uniformly. Without the pose-conditioned strategy, objects in different views may converge to inconsistent poses. As little hint is available at the start of training, the diffusion model may generate objects with different poses across view points. This inconsistency is continually introduced into NeRF training and then in return, be passed as diffusion model input to further strengthen the view difference. The diffusion model has no intent to fix such an inconsistency as these objects are individually good enough to match the text prompt and the fine-tuned object target. This problem is also known as the Janus Problem, which commonly exists in many text-to-3D NeRF generation or editing methods.

Due to the nature of the NeRF parameters, a view rendered closed to the already-trained views should contain similar but slightly blurry and distorted object content. It gives the diffusion model enough hints to generate objects of similar and more correct pose in these nearby views. The consistent information thus progressively expands to the entire dataset and is fused as 3D in NeRF.

Periodic dataset updates. We observe that periodically updating the data samples through the diffusion model is important to achieve high-quality results. In fact, when we keep the data samples from the pose-conditioned dataset updates fixed in the entire training, we found that although objects in the rendering have fairly accurate positions and orientations, their texture and geometry suffer from inconsistencies. Having periodically updates on the data samples through diffusion model fixes these minor defects and facilitates a more consistent convergence of the training dataset.

Noise strength. Unsuitable noise strength used by the diffusion model (in Eq. 6) results in various failures. Too-strong

noise disables the NeRF rendering information being passed to the diffusion model, which diminishes the effects of all dataset updating strategies and falls into the corresponding consistency problems. On the other hand, having too-low noise hinders the diffusion model from converging to plausible content. The effectiveness of each model inference is reduced, and the artifacts remain even after thousands of model inferences. We also test inferencing with random strengths within a range $[0.02, 0.98]$ as in Instruct-NeRF2NeRF [13] and found that high-noise inference dominates and leads to similar pose inconsistency. We empirically found that $\tau = 0.35$ is a noise strength that works for most of our editing samples.

5. Conclusions

In this paper, we introduce a new method for object insertion and object removal with neural radiance fields. Our method is built upon the novel idea of iterative dataset updates that leverages a text-to-image diffusion model to fuse object into background images and a pose-conditioned dataset updates strategy to stabilize the NeRF training when the object manipulation is progressively introduced into the training process. Our method shows promising results on a diverse set of test scenes and objects.

Our method is not without limitations. First, since we consider images from the fine-tuned diffusion model as the final output, we can face similar flickering problems as in the video editing results produced by diffusion models [24, 47, 65]. We leave this issue for future work, which can potentially be addressed by SOTA video translation methods [69]. Next, extending our method to support more tasks, e.g., object translation and rotation, would result in more robust object editing for NeRFs. Finally, it is of great interest to consider a theoretical formulation of dataset updates to better understand the convergence of NeRF training for editing purposes.

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. 2
- [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/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. 5
- [3] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023. 2
- [4] Ronald A Castellino. Computer aided detection (cad): an overview. *Cancer Imaging*, 5(1):17, 2005. 1

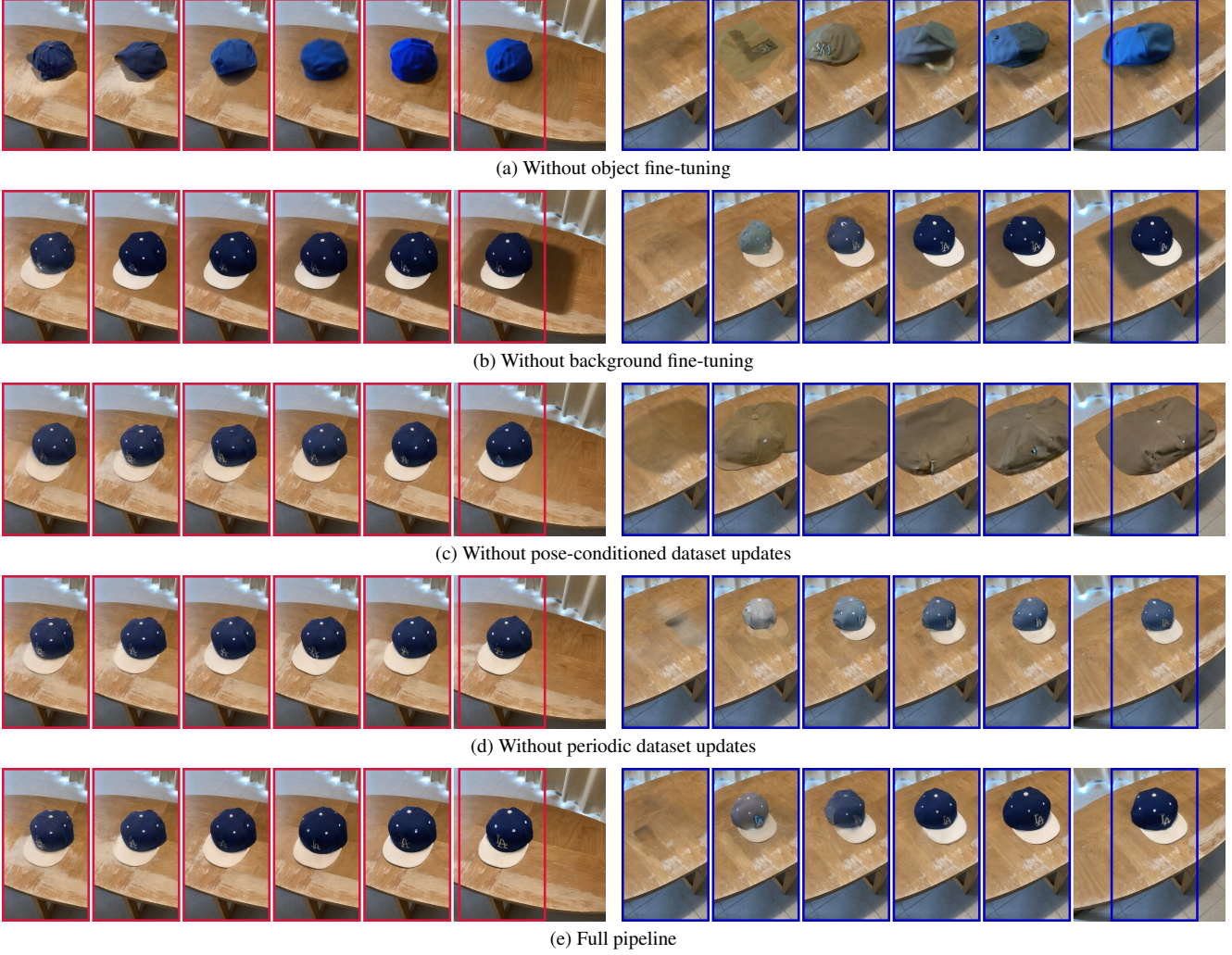


Figure 5. Ablation study that shows the effectiveness of our proposed modules. Each row corresponds to the results of our pipeline being trained without a specific module. The left and right columns of images being bounded by red and blue are the two different views of a same edited scene. For the red or blue column, from left to right are the results from increasing training steps, where the right most image is the final output. Note that the red column is a view near the starting view, which converges faster than the blue column from a farther view (except in (c) all views converge equally).

- [5] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 2
- [6] Hansheng Chen, Jiatao Gu, Anpei Chen, Wei Tian, Zhuowen Tu, Lingjie Liu, and Hao Su. Single-stage diffusion nerf: A unified approach to 3d generation and reconstruction. *arXiv preprint arXiv:2304.06714*, 2023. 6
- [7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019. 2
- [8] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 1
- [9] Epic Games. Unreal engine. 1
- [10] SM Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S Morcos, Marta Garnelo, Avraham Ruderman, Andrei A Rusu, Ivo Danihelka, Karol Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018. 2
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 1, 2
- [12] John K Haas. A history of the unity game engine. 2014. 1
- [13] Ayaan Haque, Matthew Tancik, Alexei A Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing

- 3d scenes with instructions. *arXiv preprint arXiv:2303.12789*, 2023. 1, 2, 3, 5, 6, 7, 8, 9
- [14] Jeffrey Harper. *Mastering Autodesk 3ds Max 2013*. John Wiley & Sons, 2012. 1
- [15] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. 6
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 2
- [17] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2821–2830, 2018. 2
- [18] Yi-Hua Huang, Yue He, Yu-Jie Yuan, Yu-Kun Lai, and Lin Gao. Stylizednerf: consistent 3d scene stylization as stylized nerf via 2d-3d mutual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18342–18352, 2022. 2
- [19] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 2
- [20] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *Advances in neural information processing systems*, 33:12104–12114, 2020. 2
- [21] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 2
- [22] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 2
- [23] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6007–6017, 2023. 2
- [24] Gyeongman Kim, Hajin Shim, Hyunsu Kim, Yunjey Choi, Junho Kim, and Eunho Yang. Diffusion video autoencoders: Toward temporally consistent face video editing via disentangled video encoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6091–6100, 2023. 9
- [25] Hyunsu Kim, Gayoung Lee, Yunjey Choi, Jin-Hwa Kim, and Jun-Yan Zhu. 3d-aware blending with generative nerfs. *arXiv preprint arXiv:2302.06608*, 2023. 2, 6
- [26] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. *Advances in Neural Information Processing Systems*, 35:23311–23330, 2022. 2
- [27] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012. 2
- [28] Gang Li, Heliang Zheng, Chaoyue Wang, Chang Li, Changwen Zheng, and Dacheng Tao. 3ddesigner: Towards photo-realistic 3d object generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2211.14108*, 2022. 2
- [29] Yuhao Li, Yishun Dou, Yue Shi, Yu Lei, Xuanhong Chen, Yi Zhang, Peng Zhou, and Bingbing Ni. Focaldreamer: Text-driven 3d editing via focal-fusion assembly. *arXiv preprint arXiv:2308.10608*, 2023. 2
- [30] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 300–309, 2023. 2, 3
- [31] Kai-En Lin, Yen-Chen Lin, Wei-Sheng Lai, Tsung-Yi Lin, Yi-Chang Shih, and Ravi Ramamoorthi. Vision transformer for nerf-based view synthesis from a single input image. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 806–815, 2023. 6
- [32] Miaomiao Liu, Xuming He, and Mathieu Salzmann. Geometry-aware deep network for single-image novel view synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4616–4624, 2018. 2
- [33] 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. *arXiv preprint arXiv:2303.11328*, 2023. 6, 7, 8
- [34] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. Editing conditional radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5773–5783, 2021. 2
- [35] 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/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 2
- [36] 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. 2
- [37] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Realfusion: 360deg reconstruction of any object from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8446–8455, 2023. 3
- [38] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4743–4752, 2019. 2
- [39] 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
- [40] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2
- [41] 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/CVF Conference on Computer Vision and Pattern Recognition*, pages 20669–20679, 2023. 2, 6
- [42] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 2
- [43] Thu Nguyen-Phuoc, Feng Liu, and Lei Xiao. Snerf: stylized neural implicit representations for 3d scenes. *arXiv preprint arXiv:2207.02363*, 2022. 2
- [44] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. Stylesdf: High-resolution 3d-consistent image and geometry generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13503–13513, 2022. 2
- [45] Hong-Wing Pang, Binh-Son Hua, and Sai-Kit Yeung. Locally stylized neural radiance fields. In *ICCV*, 2023. 2
- [46] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2, 5, 6
- [47] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing. *arXiv preprint arXiv:2303.09535*, 2023. 9
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2
- [49] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 1, 2
- [50] 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*, pages 10684–10695, 2022. 1, 2, 3, 6, 7, 8
- [51] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 2, 3, 4, 5, 6
- [52] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 1, 2
- [53] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 5, 7, 8
- [54] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 5, 7, 8
- [55] Katja Schwarz, Axel Sauer, Michael Niemeyer, Yiyi Liao, and Andreas Geiger. Voxgraf: Fast 3d-aware image synthesis with sparse voxel grids. *Advances in Neural Information Processing Systems*, 35:33999–34011, 2022. 2
- [56] Ivan Skorokhodov, Sergey Tulyakov, Yiqun Wang, and Peter Wonka. Epigraf: Rethinking training of 3d gans. *Advances in Neural Information Processing Systems*, 35:24487–24501, 2022. 2
- [57] Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Multi-view consistency as supervisory signal for learning shape and pose prediction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2897–2905, 2018. 2
- [58] 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*, pages 3835–3844, 2022. 2
- [59] Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Nerf-art: Text-driven neural radiance fields stylization. *IEEE Transactions on Visualization and Computer Graphics*, 2023. 2
- [60] Dongqing Wang, Tong Zhang, Alaa Abboud, and Sabine Süsstrunk. Inpaintnerf360: Text-guided 3d inpainting on unbounded neural radiance fields. *arXiv preprint arXiv:2305.15094*, 2023. 6
- [61] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. 5
- [62] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018. 2
- [63] 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. *arXiv preprint arXiv:2305.16213*, 2023. 2, 3
- [64] 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/CVF Conference on Computer Vision and Pattern Recognition*, pages 16528–16538, 2023. 2, 6
- [65] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei, Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and

- Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. *arXiv preprint arXiv:2212.11565*, 2022. 9
- [66] 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/CVF Conference on Computer Vision and Pattern Recognition*, pages 22428–22437, 2023. 2
- [67] Jiale Xu, Xintao Wang, Weihao Cheng, Yan-Pei Cao, Ying Shan, Xiaohu Qie, and Shenghua Gao. Dream3d: Zero-shot text-to-3d synthesis using 3d shape prior and text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20908–20918, 2023. 3
- [68] Yinghao Xu, Menglei Chai, Zifan Shi, Sida Peng, Ivan Skokhodov, Aliaksandr Siarohin, Ceyuan Yang, Yujun Shen, Hsin-Ying Lee, Bolei Zhou, et al. Discoscene: Spatially disentangled generative radiance fields for controllable 3d-aware scene synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4402–4412, 2023. 2
- [69] Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change Loy. Rerender a video: Zero-shot text-guided video-to-video translation. *arXiv preprint arXiv:2306.07954*, 2023. 9
- [70] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 2, 6
- [71] 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
- [72] Kai Zhang, Gernot Riegler, Noah Snaveley, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2
- [73] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 2
- [74] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 286–301. Springer, 2016. 2
- [75] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 2
- [76] Jingyu Zhuang, Chen Wang, Lingjie Liu, Liang Lin, and Guanbin Li. Dreameditor: Text-driven 3d scene editing with neural fields. *arXiv preprint arXiv:2306.13455*, 2023. 2, 6