

# Instruct-NeRF2NeRF: Editing 3D Scenes with Instructions

Ayaan Haque, Matthew Tancik, Alexei A. Efros, Aleksander Holynski, and Angjoo Kanazawa

UC Berkeley

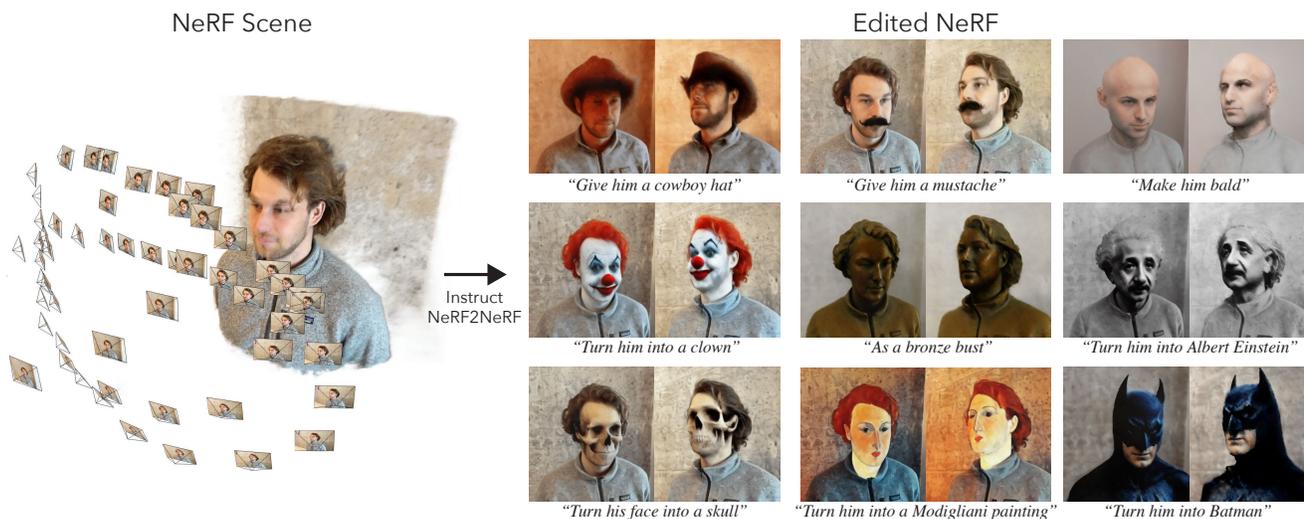


Figure 1: **Editing 3D scenes with Instructions.** We propose Instruct-NeRF2NeRF, a method for consistent 3D editing of a NeRF scene using text-based instructions. Our method can accomplish a diverse collection of local and global scene edits.

## Abstract

We propose a method for editing NeRF scenes with text-instructions. Given a NeRF of a scene and the collection of images used to reconstruct it, our method uses an image-conditioned diffusion model (InstructPix2Pix) to iteratively edit the input images while optimizing the underlying scene, resulting in an optimized 3D scene that respects the edit instruction. We demonstrate that our proposed method is able to edit large-scale, real-world scenes, and is able to accomplish more realistic, targeted edits than prior work. Result videos can be found on the project website: <https://instruct-nerf2nerf.github.io>.

## 1. Introduction

With the emergence of efficient neural 3D reconstruction techniques, capturing a realistic digital representation of a real-world 3D scene has never been easier. The process is simple: capture a collection of images of a scene

from varying viewpoints, reconstruct their camera parameters, and use the posed images to optimize a Neural Radiance Field [26]. Due to its ease of use, we expect captured 3D content to gradually replace the traditional processes of manually-generated assets. Unfortunately, while the pipelines for turning a real scene into a 3D representation are relatively mature and accessible, many of the other necessary tools for the creation of 3D assets (e.g., those needed for *editing* 3D scenes) remain underdeveloped.

Traditional processes for editing 3D models involve specialized tools and years of training in order to manually sculpt, extrude, and re-texture a given object. This process is made even more involved with the advent of neural representations, which often do not have explicit surfaces. This further motivates the need for 3D editing approaches designed for the modern era of 3D representations, particularly approaches that are similarly as accessible as the capture techniques themselves.

To this end, we propose Instruct-NeRF2NeRF, a method for editing 3D NeRF scenes that requires as input only a

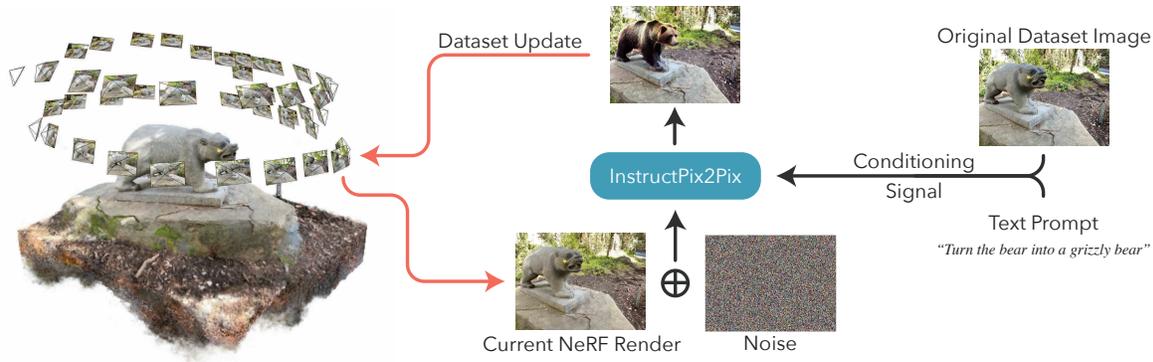


Figure 2: **Overview:** Our method gradually updates a reconstructed NeRF scene by iteratively updating the dataset images while training the NeRF: (1) an image is rendered from the scene at a training viewpoint, (2) it is edited by InstructPix2Pix given a global text instruction, (3) the training dataset image is replaced with the edited image, and (4) the NeRF continues training as usual.

text instruction. Our approach operates on a pre-captured 3D scene and ensures that the resulting edits are reflected in a 3D-consistent manner. For example, given a 3D scene capture of a person shown in Figure 1 (left), we can enable a wide variety of edits using flexible and expressive textual instruction such as “Give him a cowboy hat” or “Turn him into Albert Einstein”. Our approach makes 3D scene editing accessible and intuitive for everyday users.

Though there exist 3D generative models, the data-sources required for training these models at scale are still limited. Therefore, we instead choose to extract shape and appearance priors from a 2D diffusion model. Specifically, we employ a recent image-conditioned diffusion model, InstructPix2Pix [2], which enables instruction-based 2D image editing. Unfortunately, applying this model on individual images rendered from a reconstructed NeRF produces inconsistent edits across viewpoints. As a solution to this, we devise a simple approach similar to recent 3D generation solutions like DreamFusion [33]. Our underlying method, which we refer to as Iterative Dataset Update (Iterative DU), alternates between editing the “dataset” of NeRF input images, and updating the underlying 3D representation to incorporate the edited images.

We evaluate our approach on a variety of captured NeRF scenes, validating our design choices by comparing with ablated variants of our method, as well as naïve implementations of the score distillation sampling (SDS) loss proposed in DreamFusion [33]. We also qualitatively compare our approach to a concurrent text-based stylization approach [47]. We demonstrate that our method can accomplish a wide variety of edits on people, objects, and large-scale scenes.

## 2. Related Work

**Physical Editing of NeRFs** NeRFs [26] are a popular approach for generating photorealistic novel views of a scene captured by calibrated photographs and have been extended

in many follow-up works [42]. However, editing NeRFs remains a challenge due to their underlying representation. One approach is to impose physics-based inductive biases in its optimization process to enable changes in materials or scene lighting [44, 1, 40, 27, 25]. Alternatively, one can specify bounding boxes [30, 50], to allow easy compositing of different objects [52] as well as spatial manipulations and geometry deformations [51]. A recent work, ClimateNeRF [18], extracts rough geometry from a NeRF and uses physical simulation to apply weather changes such as snowing and flooding. Most physically-based edits revolve around changing physical properties of the reconstructed scene, or performing physical simulation. In this work, we instead focus on enabling arbitrary creative edits.

**Artistic Stylization of NeRFs** Following the literature from image stylization [8, 7], recent works have explored artistic 3D stylization of NeRFs [5, 12, 13, 28, 53, 49]. While these approaches can obtain 3D-consistent stylizations of a scene, they primarily focus on global scene appearance changes and usually require a reference image. Other works have explored the use of latent representations from visual language models such as CLIP [34]. EditNeRF [20] explores editing NeRFs by manipulating latent codes learned from object categories in a synthetic dataset. To increase usability (and as explored in other 3D domains such as meshes [24, 10]), ClipNeRF [46] and NeRF-Art [47] extend this line of work by encouraging similarity between CLIP embeddings of the scene and a short text prompt. A limitation of these CLIP-based approaches is their inability to incorporate localized edits. Methods such as Distilled Feature Fields [15] and Neural Feature Fusion Fields [43] distill 2D features from pre-trained models such as LSeg [17] and DINO [4] into the radiance fields, which enable specification of regions. These approaches allow for localized CLIP-guided edits, 3D spatial transformations [46], or localized scene removal [43] specified either

by language or a reference image. In this work, we offer a complementary approach to editing 3D scenes based on intuitive, purely language-based editing instructions. While masking enables specific local changes, instructional edits provide intuitive high-level instructions that can make more flexible and holistic changes to the appearance or geometry of a single object or the entire scene. We enable mask-free instructional edits by taking advantage of recent instruction-based 2D image-conditioned diffusion model [2], resulting in a purely-language-based interface that enables a wider range of intuitive and content-aware 3D editing.

**Generating 3D Content** Recent progress in pre-trained large-scale models has enabled rapid progress in the domain of generating 3D content from scratch, either by optimizing radiance fields through vision-language models like CLIP [14, 16] or via text-conditioned diffusion models [35, 37, 36] as presented in DreamFusion [33] and its follow-ups [48, 19, 23]. While these approaches can generate 3D models from arbitrary text prompts, they lack (1) fine-grained control over the synthesized outputs, (2) the ability to generalize to scenes (*i.e.*, anything beyond a single object isolated in space), and (3) any grounding in reality, producing entirely synthesized creations. Concurrent works such as RealFusion [21] and SparseFusion [55] explore grounding by providing one or few input images, where the unseen parts are hallucinated. In all of these approaches, a central challenge is congealing the inconsistent outputs of a 2D diffusion model into a consistent 3D scene. In this work, instead of creating new content, we focus on editing *real* captured NeRFs of *fully observed scenes* using 2D diffusion priors. One advantage of editing an existing NeRF scene (as opposed to generating 3D content from scratch) is that the captured images are by definition 3D consistent, suggesting that generated imagery should naturally be more consistent. This also helps avoid certain design decisions that result in the cartoon-ish appearance commonly seen in unconditional 3D content generation methods [33, 48, 19].

**Instruction as an Editing Interface** With the rise of large-language models (LLMs) like GPT [3] and ChatGPT [29], natural language is emerging as the next “programming language” for specifying complex tasks. LLMs allow for the abstraction of a series of low-level specifications into an intuitive and user-friendly interface through the use of language, specifically *instructions* [31]. InstructPix2Pix [2] demonstrates the effectiveness of instructions in 2D image tasks, as do other works in other domains such as robotic navigation [11]. We propose the first work that demonstrates instructional guidance in the realm of 3D editing. This is particularly significant given the difficulty of the task, which has typically required specialized tools and years of experience. By using natural language instructions,

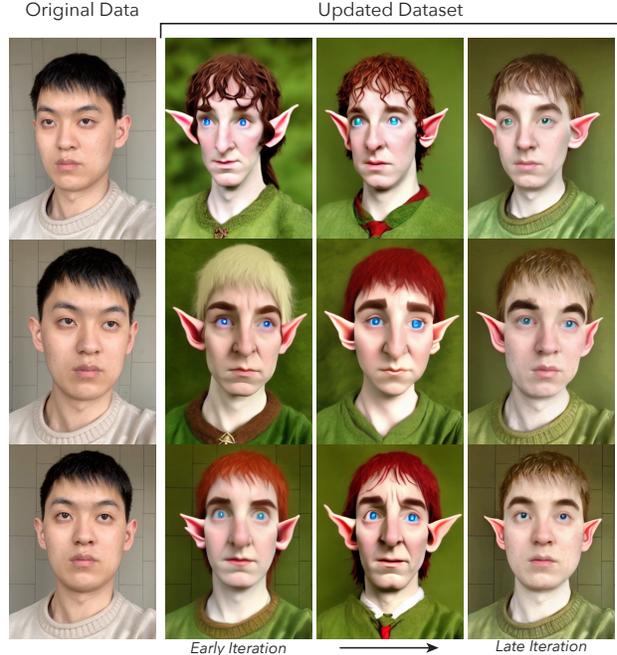


Figure 3: **Dataset Evolution:** At the start of training, the edited images perform the requested edit, but are often inconsistent. After iteratively training the NeRF and updating the training dataset, the images gradually become more 3D consistent.

even novice users can achieve high-quality results without additional tools or specialized knowledge.

### 3. Method

Our method takes as input a reconstructed NeRF scene along with its corresponding source data: a set of captured images, their corresponding camera poses, and camera calibration (typically from a structure-from-motion system, such as COLMAP [38]). Additionally, our method takes as input a natural-language editing instruction, e.g., “*turn him into Albert Einstein*”. As output, our method produces an edited version of the NeRF subject to the provided edit instruction, as well as edited versions of the input images.

Our method accomplishes this task by iteratively updating the image content at the captured viewpoints with the help of a diffusion model, and subsequently consolidating these edits in 3D through standard NeRF training. Our work builds off recent advances in diffusion models for image editing, specifically InstructPix2Pix [2], which proposes an image-and-text conditioned diffusion model trained to edit natural images using human-provided instructions.

#### 3.1. Background

**Neural radiance fields** Neural radiance fields (NeRFs) [26] are a compact and convenient representation for reconstructing and rendering a volumetric 3D

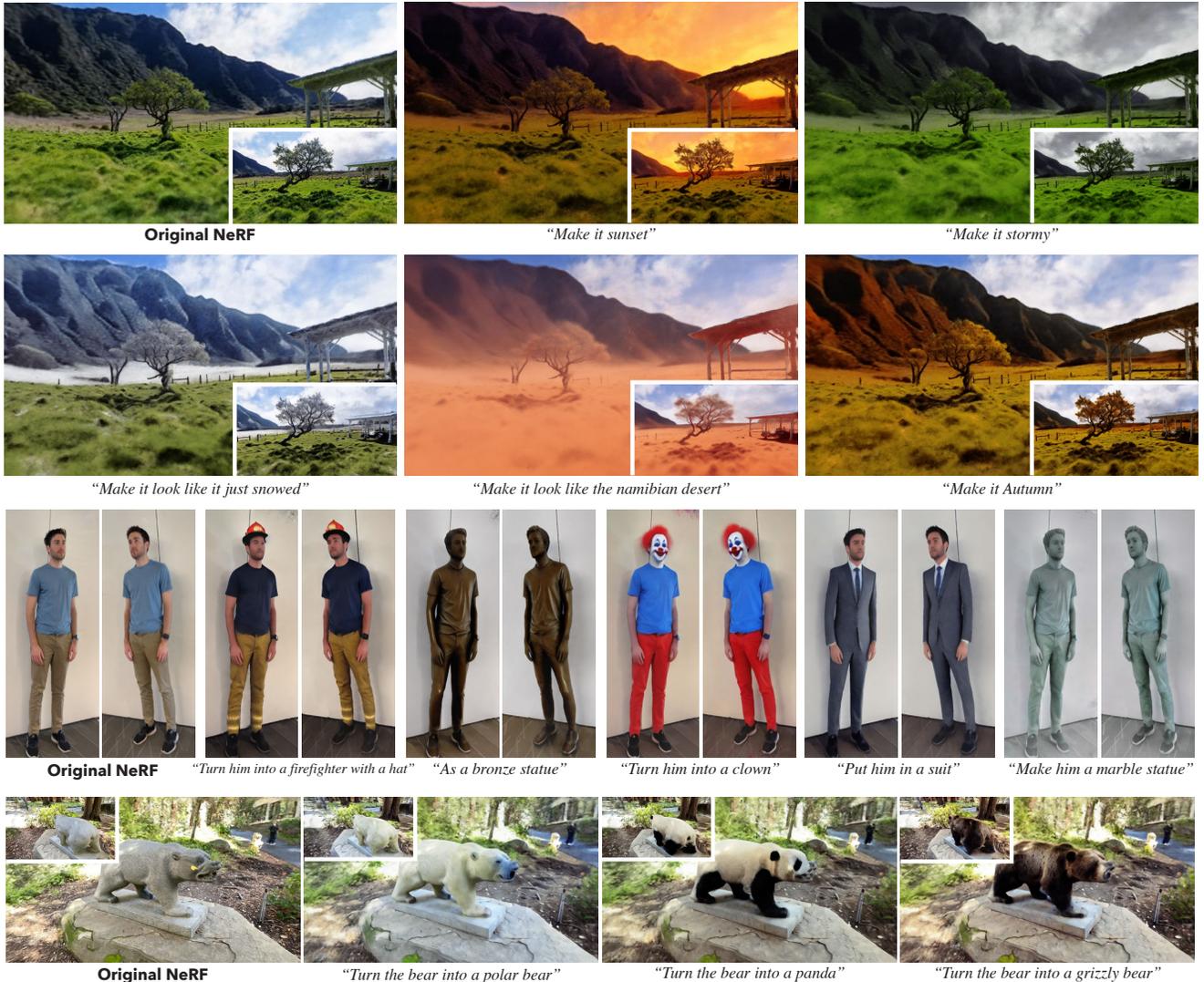


Figure 4: **Qualitative Results:** Our method is able to perform a variety of diverse contextual edits on real scenes, including environmental changes, like adjusting the time of day, and even more localized changes that modify only a specific object in the scene.

scene. A NeRF is parameterized by 3D positions  $(x, y, z)$  and viewing directions  $(\theta, \phi)$  for samples in a field. Each sample is processed to produce a color and density  $(c, \sigma)$ , which can be composited along a ray to produce a 2D pixel color. A NeRF is optimized using a collection of captured images and their corresponding camera parameters, which include both calibration and extrinsic pose/orientation. These camera parameters can be used to extract a per-pixel world-space ray parameterization that describes the 3D center  $\mathbf{o}$  and direction  $\mathbf{d}$  of the camera ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  that corresponds to each pixel in each image. These rays with their associated pixel colors are used to optimize the NeRF. The typical process of training a NeRF [26] involves selecting a subset of rays  $\mathbf{r}$ , rendering the NeRF’s current estimate of the color along this ray  $\hat{C}(\mathbf{r})$ , and computing

a loss relative to captured pixel color  $\mathcal{L}(C(\mathbf{r}), \hat{C}(\mathbf{r}))$ . In practice, in the interest of reliable optimization, rays are selected at random from a variety of viewpoints, to ensure the 3D positions of reconstructed scene objects are sufficiently well-constrained. To render a novel viewpoint, a collection of rays are sampled corresponding to all the pixels in that novel image, and the resulting color values  $\hat{C}(\mathbf{r})$  are arranged into a 2D frame.

**InstructPix2Pix** Denoising diffusion models [39, 9] are generative models that learn to gradually transform a noisy sample towards a modeled data distribution. InstructPix2Pix [2] is a diffusion-based method specialized for image editing. Conditioned on an RGB image  $c_I$  and



Figure 5: **Guidance Scale:** By varying the image guidance we can control how much the edit looks like the original scene. Note that these are renderings from the edited 3D scenes.

a text-based editing instruction  $c_T$ , and taking as input a noised image (or pure noise)  $z_t$ , the model aims to produce an estimate of the edited image  $z_0$  (an edited version of  $c_I$  subject to the instruction  $c_T$ ). Formally, the diffusion model predicts the amount of noise present in the input image  $z_t$ , using the denoising U-Net  $\epsilon_\theta$  as:

$$\hat{\epsilon} = \epsilon_\theta(z_t; t, c_I, c_T) \quad (1)$$

This noise prediction  $\hat{\epsilon}$  can be used to derive  $\hat{z}_0$ , the estimate of the edited image. This denoising process can be queried with a noisy image  $z_t$  at any timestep  $t \in [0, T]$ , i.e., containing any amount of noise, up to a pure noise image  $z_T$ . Larger amounts of noise, i.e., larger values of  $t$ , will produce estimates of  $\hat{z}_0$  with more variance, whereas smaller  $t$  values will produce lower variance estimates with more adherence to the visible image signal in  $z_t$ .

In practice, InstructPix2Pix is based on a latent diffusion model [36], i.e., the diffusion process operates entirely on an encoded latent domain. This means that the above-defined variables  $c_I, z_0$  are all latent images created by encoding an RGB image, i.e.,  $\mathcal{E}(I)$ . Similarly, to produce an RGB image from the diffusion model, one must also decode the predicted  $\hat{z}_0$  latents via the decoder  $\hat{I} = \mathcal{D}(\hat{z}_0)$ .

### 3.2. Instruct-NeRF2NeRF

Given a reconstructed NeRF scene (including the corresponding dataset of calibrated images), as well as a text instruction, we fine-tune the reconstructed model towards an edit instruction to produce an edited version of that NeRF. An overview is provided in Fig. 2.

Our method works through an alternating update scheme, in which the training dataset images are iteratively updated using a diffusion model and are subsequently consolidated into the globally consistent 3D representation by training the NeRF on these updated images. This iterative process allows for gradual percolation of the diffusion priors into the 3D scene. Although this process can enable significant edits to the scene, our use of an image-conditioned diffusion model (InstructPix2Pix) helps in maintaining the structure and identity of the original scene.

In this section, we first describe our use of Instruct-Pix2Pix in the process of editing a single dataset image,

then describe our iterative procedure for gradually updating dataset images and refining the reconstructed NeRF.

**Editing a rendered image** We use InstructPix2Pix [2] to edit each dataset image. It takes three inputs: (1) an input conditioning image  $c_I$ , a text instruction  $c_T$ , and a noisy input  $z_t$ . To update a dataset image at viewpoint  $v$ , we use the unedited image  $I_0^v$  for  $c_I$ , which will typically be the originally captured image at this viewpoint, or if the viewpoint was not captured physically, a render from the NeRF before any edits were made. For  $z_t$ , as in SDEdit [22], we input a noised version of the current render at optimization step  $i$ , i.e., a linear combination of  $\mathcal{N}(0, 1)$  and  $z_0 = \mathcal{E}(I_i^v)$ . For simplicity, we denote the process of replacing an image  $I_i^v$  as  $I_{i+1}^v \leftarrow U_\theta(I_i^v, t; I_0^v, c_T)$ , where a noise level  $t$  is chosen at random from a constant range  $[t_{\min}, t_{\max}]$ . We define  $U_\theta$  as the DDIM sampling process, with a fixed number of intermediate steps  $s$  taken between initial timestep  $t$  and 0.

This process mixes two sources of information: the diffusion model aims to edit the original image  $I_0^v$  according to the instruction  $c_T$ , while the noised image passed to the diffusion U-Net  $z_t$  is only partially noised (with some  $t < T$ ), such that the rendering of the current global 3D model influences the diffusion model’s final estimate for  $z_0$  (the image which will replace the dataset image at viewpoint  $v$ ). A key thing to note is that while our method continually repeats the process of rendering from the NeRF, editing the image, and updating the NeRF, the diffusion model is conditioned on the *un-edited* images, and thus remains grounded, preventing the characteristic drift commonly associated with recurrent synthesis.

**Iterative Dataset Update** The core component of our method is an alternating process through which images are rendered from the NeRF, updated by the diffusion model, and subsequently used to supervise the NeRF reconstruction. We refer to this process as the Iterative Dataset Update (*Iterative DU*).

When optimization begins, our image dataset consists of the originally captured images from a range of viewpoints denoted as  $v$ , which we represent as  $I_0^v$ . These images are cached separately and used as conditioning for the diffusion model at all stages. At each iteration, we perform a number of image updates  $d$ , followed by a number of NeRF updates  $n$ . Image updates are performed sequentially in a random ordering of  $v$  determined at the start of training. NeRF updates always sample a set of random rays from the entire training dataset, such that the supervision signal contains a mixture of *old* information and pixels from recently updated dataset images.

Our editing process results in sudden replacement of dataset images with their edited counterpart. At early iterations, these images may perform inconsistent edits (as

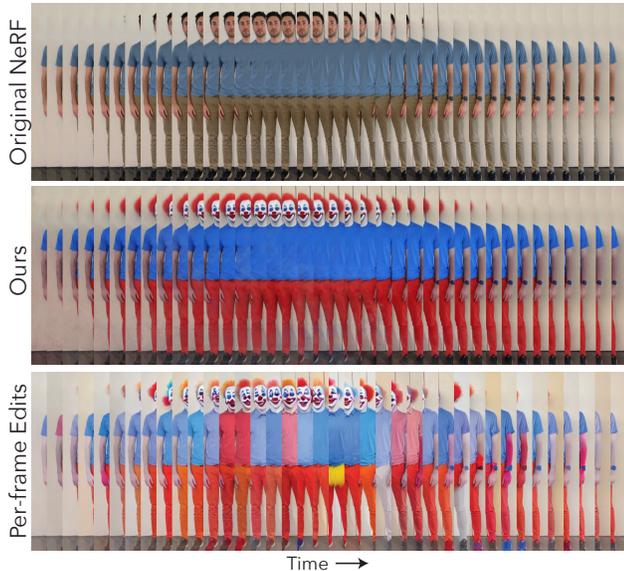


Figure 6: **Consistency:** Vertical slices of a rendered novel camera path show the consistency across varying viewpoints. The original NeRF rendering (top) is quite consistent, similar to our edited result (middle), using the prompt “turn him into a clown”. Conversely, running InstructPix2Pix [2] on each rendered frame independently results in notable inconsistency, such as varying hair and shirt colors.

InstructPix2Pix does not typically perform consistent edits across different viewpoints). Over time, as images are used to update the NeRF and progressively re-rendered and updated, they begin to converge on a globally consistent depiction of the edited scene. Examples of this evolution process can be seen in Figure 3.

This process is similar to the approach proposed in SNeRF [28], where the images are updated through style transfer in every alternate iteration. Unlike SNeRF, our iterative DU retains edited images across NeRF updates, effectively performing semi-permanent updates to the training dataset. This process can also be interpreted as a variant of the score distillation sampling (SDS) loss from DreamFusion [33], where instead of updating a discrete set of images at each step, each gradient update contains a random mixture of rays distributed across many viewpoints, and the computed gradients along these rays may not be from the most recent NeRF state. The use of iterative DU is aimed at maximizing the diversity of training ray viewpoints in each iteration, a choice that we find greatly improves both training stability and efficiency. In the following section, we provide a comparison to a naïve adaptation of the SDS loss to our application.

### 3.3. Implementation details

As the underlying NeRF implementation, we use the ‘nerfacto’ model from NeRFStudio [41]. The strength and consistency of the updates performed by the diffusion model are determined by several parameters. Among these are the values for  $[t_{\min}, t_{\max}] = [0.02, 0.98]$ , which define the amount of noise (and therefore the amount signal retained from the original images). Regardless of  $t$ , we always sample our denoised image with 20 denoising steps. The diffusion model has additional parameters, such as the classifier-free guidance weights corresponding to the text and image conditioning signals. For these, we can use the default values of  $s_I = 1.5$  and  $s_T = 7.5$ , or offer the user the ability to hand-tune the guidance weight on an image to achieve the optimal edit strength before performing our NeRF optimization process. The results shown in the paper use manually selected guidance values, but adjusting these can result in varying degrees of scene edits, as shown in Figure 5. All other 3D hyperparameters are fixed for all experiments. During optimization, for the sake of efficiency, we update one image at a time, i.e.,  $d = 1$  and  $n = 10$ . For NeRF training, we use L1 and LPIPS [54] losses. We train our method for a maximum of 30k iterations, which takes roughly an hour on a single NVIDIA Titan RTX (15GB of memory). More details are provided in the appendix.

## 4. Results

We conduct experiments on real scenes optimized using Nerfstudio [41]. We edit a variety of scenes that vary in complexity: 360 scenes of environments and objects, faces, and full-body portraits. The scenes were captured using both a smartphone and a mirrorless camera. The camera poses were extracted using either COLMAP [38] or through the PolyCam [32] app. The size of each dataset ranges from 50-300 images. First, we evaluate our approach through a variety of qualitative evaluations. To validate our design choices, we compare against a set of ablation baselines both qualitatively and quantitatively. Additionally, we provide visual comparisons to concurrent work NeRF-Art [47].

### 4.1. Qualitative Evaluation

**Editing 3D Scenes** Our qualitative results are shown in Figure 1 and Figure 4. For each edit, we show multiple views to illustrate the 3D consistency. On the portrait capture in Figure 1, we are able to achieve a broad range of edits varying from global (“Turn him into a Modigliani painting”) to locally specific edits (“Turn his face into a skull”). Although adding a completely new object is as challenging as the task of DreamFusion, our approach is able to add contextual elements such as “Give him a cowboy hat” and “mustache”. Moreover, our method is able to dress the person to some degree, such as those illustrated on the full-



Figure 7: **Comparison with NeRF-Art:** We compare with CLIP-based method NeRF-Art on sequences and edits provided in their paper.

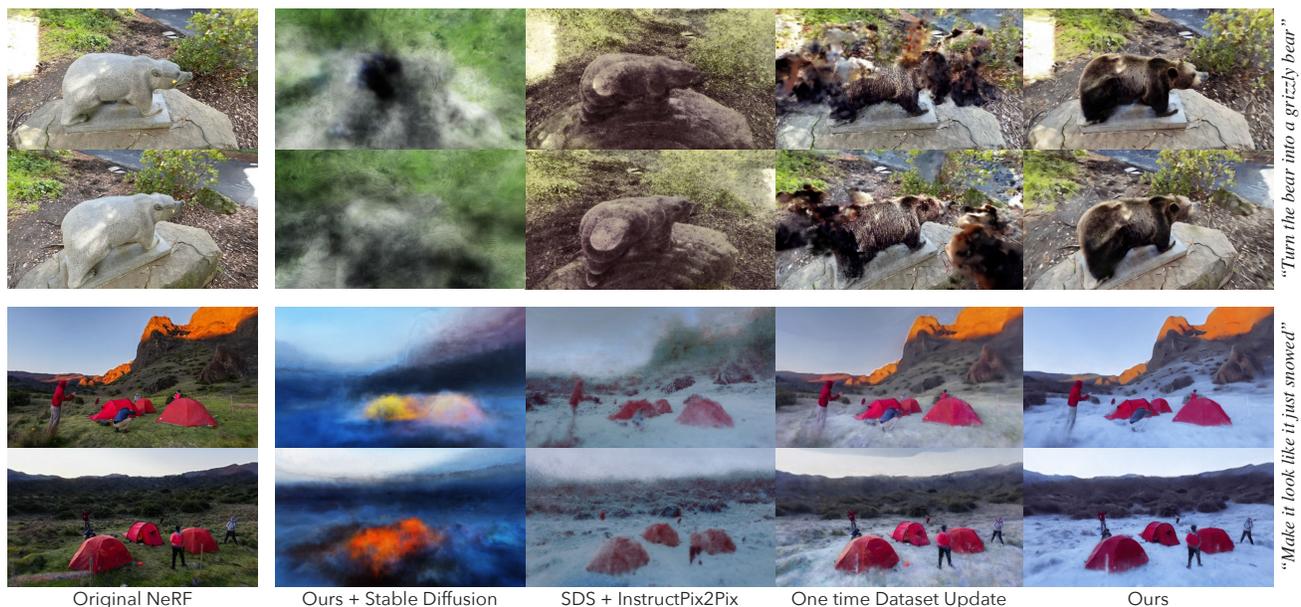


Figure 8: **Baseline Comparisons:** We compare our model with a collection of variants described in Section 4.1.

body portrait in Figure 4, third row. It can achieve material changes such as “As a bronze bust” and “Make him a marble statue”. In the “bronze” cases a subtle amount of view-dependent changes are also captured. Our approach is also able to turn portraits into notable figures such as Einstein and fictional characters like “Batman”. These edits also extend to subjects other than people, like changing a bear statue into a real polar bear, panda, and grizzly bear (Figure 4, last row). Most notably, these edits also apply to large-scale scenes (Figure 4, first row, Figure 8, bottom), and support instructions that modify the time of the day, seasons, and other conditions such as snow and desert.

**Ablation Study** We validate our design choices by comparing our approach to the following variants. The qualitative differences are shown in Figure 8:

*Per-frame Edit.* Our most naive baseline is to apply Instruct-Pix2Pix [2] independently on every frame of a novel path rendered by the original NeRF. We use the rendered images as  $c_I$ , and the same text instruction as  $c_T$ . For  $z_t$ , we use pure noise. Despite the fact that the conditioning images are 3D consistent, the resulting edited images have significant variance that is inconsistent across different views. We illustrate this inconsistency in Figure 6, where we pan a camera across the scene and concatenate a slice from each frame to create an image. See supplemental video for examples.

*One time Dataset Update.* In this baseline, we perform a single Dataset Update step, in which all training images are edited once, and the NeRF is trained until convergence on those edited images. The quality of this baseline depends largely on the 3D consistency of the per-frame editing results. While this approach can sometimes yield decent results, in a majority of cases, the initial edited 2D images are

	CLIP Text-Image Direction Similarity $\uparrow$	CLIP Direction Consistency $\uparrow$
Per-frame IP2P [2]	<b>0.1603</b>	0.8185
One-time DU	0.1157	0.8823
SDS w/ IP2P [2]	0.0266	<i>0.9160</i>
Ours	<i>0.1600</i>	<b>0.9191</b>

Table 1: **Quantitative Evaluation.** Although edits are subjective, we provide quantitative metrics that evaluate the alignment of the edits to the text and consistency between subsequent frames in the CLIP space. Our approach results in similar CLIP similarity as per-frame edit, while achieving best consistency in CLIP space.

largely inconsistent, leading to blurry and artifact-filled 3D scenes, as shown in Figure 8. This problem is even more prominent when contextual objects are added to the portraits.

*DreamFusion (text-conditioned diffusion).* The next approach is to naively apply DreamFusion optimization to an existing NeRF scene. Specifically, starting from a NeRF initialized by the density and appearance obtained from real images, we apply SDS loss [33] using StableDiffusion [36], a text-only diffusion model. We observed that this method quickly diverges and thus we do not include qualitative results in the paper. The reason for this divergence is that, in this setting, every image needs a textual description of the scene, and it becomes difficult to find an exact textual description that matches a scene across all views, especially for those with 360-degree coverage. This experiment highlights the importance of image conditioning.

*SDS + InstructPix2Pix.* If instead, we use an image-conditioned generative model, InstructPix2Pix, with the SDS loss from the previous variant, we are able to circumvent the requirements for an accurate text description of the whole scene. Unlike the text-conditioned variant, this approach does not diverge, but results in a 3D scene with more artifacts, as seen in Figure 8, third column. We largely attribute this to the fact that the standard SDS samples rays from a small collection of full images, which makes optimization more unreliable than sampling rays randomly across all viewpoints.

*Ours + StableDiffusion.* Finally, we compare our approach (with Iterative DU), but using StableDiffusion instead of InstructPix2Pix. This approach suffers from similar issues as seen in the DreamFusion baseline, because of the lack of image conditioning. Although it doesn’t diverge, the resulting scene is blurry, and the 3D density is not coherent. Qualitative results can be seen in Figure 8, first column.

**Comparisons with NeRF-Art.** We provide a qualitative comparison against concurrent work NeRF-Art [47]. Although their training code is unavailable, we use their provided custom-captured scenes and perform similar edits us-



Figure 9: **Limitations:** InstructPix2Pix cannot always perform the desired edit (top), and thus our method does not perform an edit. Sometimes InstructPix2Pix produces correct, but inconsistent edits in 2D that our method fails to consolidate in 3D (bottom).

ing our method. A comparison of their provided scene is shown in Figure 7. Note that their text inputs are not instructions, leaving the model with ambiguity on what exactly to edit. For instance, in their example of “Van Gogh”, it’s unclear whether the model should create a painting in the style of Van Gogh or make the face look like Van Gogh’s face. Since edits are subjective, we leave it to the readers to determine their preference for these edits and provide this as a reference to a competitive state-of-the-art.

## 4.2. Quantitative Evaluation

Editing is fundamentally a subjective task. Thus, we mostly rely on various types of qualitative evaluation. We recommend the reader to evaluate the performance through the supplemental videos. Nevertheless, inspired by the evaluation protocols in InstructPix2Pix, we report auxiliary quantitative metrics over 10 total edits across two scenes, measuring (1) the alignment of the performed 3D edit with the text instruction (as shown in InstructPix2Pix and StyleGAN-Nada [6]) and (2) the temporal consistency of the performed edit across views, shown in Table 1. The latter is a novel metric, similar to the CLIP directional similarity, but measuring the directional similarity between pairs of original and edited images in adjacent frames of novel rendered camera paths. More details are provided in the appendix.

## 4.3. Limitations

Our method inherits many of the limitations of InstructPix2Pix, such as the inability to perform large spatial manipulations. Furthermore, as in DreamFusion, our method uses a diffusion model on a single view at a time, and thus may suffer from similar artifacts, such as double faces on added objects. We demonstrate examples of two types of failure cases in Figure 9: (1) InstructPix2Pix fails to perform the edit in 2D, and therefore our method

fails in 3D, and (2) InstructPix2Pix succeeds at editing in 2D, but has large inconsistencies that our method fails to consolidate in 3D.

Specifically, our method is effective at instruction-driven, contextual, large-scale edits, which includes 1) editing textures, 2) replacing objects, and 3) changing global properties of a scene, among others. However, adding entirely new objects to the scene (such as adding a cup on a table) is challenging for various reasons. InstructPix2Pix struggles to add content into empty regions, and even when it is successful, it often places objects in different locations in different images, making 3D reconstruction a challenge. Similarly, InstructPix2Pix struggles to remove objects without replacing them with similarly salient (but also view-inconsistent) content, resulting in similar artifacts. We contend that as diffusion models improve at adding content, removing content, and overall image manipulation, our Iterative DU framework for NeRF editing will similarly improve. Limitations are further discussed in the appendix.

## 5. Conclusion

In this paper, we have introduced Instruct-NeRF2NeRF, a promising step towards the democratization of 3D scene editing for everyday users. Our method enables intuitive and accessible NeRF scene editing using natural text instructions. We operate on pre-captured NeRF scenes, ensuring that any resulting edits maintain 3D-consistency. We showed our method’s results on a variety of captured NeRF scenes and demonstrated its ability to accomplish a wide range of edits on people, objects, and large-scale scenes.

## 6. Acknowledgements

We thank our colleagues for their insightful feedback helpful discussions, in particular Ethan Weber, Frederik Warburg, Ben Poole, Richard Szeliski, Jon Barron, Alexander Kristoffersen, Rohan Mathur, Alejandro Escontrela, and the Nerfstudio team.

## References

- [1] Mark Boss, Varun Jampani, Raphael Braun, Ce Liu, Jonathan T. Barron, and Hendrik P.A. Lensch. Neural-pil: Neural pre-integrated lighting for reflectance decomposition. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [2] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions. In *CVPR*, 2023.
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- [5] Pei-Ze Chiang, Meng-Shiun Tsai, Hung-Yu Tseng, Weisheng Lai, and Wei-Chen Chiu. Stylizing 3d scene via implicit representation and hypernetwork, 2021.
- [6] Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *arXiv preprint arXiv:2108.00946*, 2021.
- [7] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016.
- [8] Aaron Hertzmann. Painterly rendering with curved brush strokes of multiple sizes. In *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, pages 453–460, 1998.
- [9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [10] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *ACM Transactions on Graphics (TOG)*, 41(4):1–19, 2022.
- [11] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, London, UK, 2023.
- [12] Hsin-Ping Huang, Hung-Yu Tseng, Saurabh Saini, Maneesh Singh, and Ming-Hsuan Yang. Learning to stylize novel views. *arXiv preprint arXiv:2105.13509*, 2021.
- [13] 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.
- [14] Ajay Jain, Ben Mildenhall, Jonathan T. Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. *CVPR*, 2022.
- [15] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. In *Advances in Neural Information Processing Systems*, volume 35, 2022.
- [16] Han-Hung Lee and Angel X Chang. Understanding pure clip guidance for voxel grid nerf models. *arXiv preprint arXiv:2209.15172*, 2022.
- [17] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022.
- [18] Yuan Li, Zhi-Hao Lin, David Forsyth, Jia-Bin Huang, and Shenlong Wang. Climatenerf: Physically-based neural rendering for extreme climate synthesis. *arXiv e-prints*, pages arXiv–2211, 2022.

- [19] 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. *arXiv preprint arXiv:2211.10440*, 2022.
- [20] 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.
- [21] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Realfusion: 360° reconstruction of any object from a single image. *arXiv e-prints*, pages arXiv–2302, 2023.
- [22] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [23] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. *arXiv preprint arXiv:2211.07600*, 2022.
- [24] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. *arXiv preprint arXiv:2112.03221*, 2021.
- [25] Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul P. Srinivasan, and Jonathan T. Barron. NeRF in the dark: High dynamic range view synthesis from noisy raw images. *CVPR*, 2022.
- [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.
- [27] Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Müller, and Sanja Fidler. Extracting triangular 3d models, materials, and lighting from images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8280–8290, 2022.
- [28] Thu Nguyen-Phuoc, Feng Liu, and Lei Xiao. Snerf: stylized neural implicit representations for 3d scenes. *ACM Transactions on Graphics (TOG)*, 41(4):1–11, 2022.
- [29] OpenAI. ChatGPT.
- [30] Julian Ost, Fahim Mannan, Nils Thuerey, Julian Knodt, and Felix Heide. Neural scene graphs for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2856–2865, 2021.
- [31] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [32] Polycam. Polycam - lidar & 3d scanner for iphone & android.
- [33] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv*, 2022.
- [34] 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.
- [35] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [36] 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.
- [37] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- [38] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.
- [39] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.
- [40] Pratul P. Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T. Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *CVPR*, 2021.
- [41] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, David McAllister, and Angjoo Kanazawa. Nerfstudio: A modular framework for neural radiance field development. *arXiv preprint arXiv:2302.04264*, 2023.
- [42] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Wang Yifan, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, et al. Advances in neural rendering. In *Computer Graphics Forum*. Wiley Online Library, 2022.
- [43] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural Feature Fusion Fields: 3D distillation of self-supervised 2D image representations. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2022.
- [44] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T. Barron, and Pratul P. Srinivasan. Ref-NeRF: Structured view-dependent appearance for neural radiance fields. *CVPR*, 2022.
- [45] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>, 2022.

- [46] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. *arXiv preprint arXiv:2112.05139*, 2021.
- [47] Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Nerf-art: Text-driven neural radiance fields stylization. *arXiv preprint arXiv:2212.08070*, 2022.
- [48] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. *arXiv preprint arXiv:2212.00774*, 2022.
- [49] Qiling Wu, Jianchao Tan, and Kun Xu. Palettenerf: Palette-based color editing for nerfs. *arXiv preprint arXiv:2212.12871*, 2022.
- [50] Hong-Xing Yu, Leonidas J Guibas, and Jiajun Wu. Unsupervised discovery of object radiance fields. *arXiv preprint arXiv:2107.07905*, 2021.
- [51] 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.
- [52] Jiakai Zhang, Xinhang Liu, Xinyi Ye, Fuqiang Zhao, Yanshun Zhang, Minye Wu, Yingliang Zhang, Lan Xu, and Jingyi Yu. Editable free-viewpoint video using a layered neural representation. *ACM Transactions on Graphics (TOG)*, 40(4):1–18, 2021.
- [53] Kai Zhang, Nick Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snavely. Arf: Artistic radiance fields, 2022.
- [54] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.
- [55] Zhizhuo Zhou and Shubham Tulsiani. Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction, 2022.

## A. Additional implementation details

The primary input to our method is a NeRF reconstruction of a real scene. This reconstruction is obtained using the ‘nerfacto’ model from NeRFStudio [41], trained for 30,000 iterations per scene. Before running InstructNeRF2NeRF, we re-initialize the optimizers in NeRFStudio. In order to specify edits, we use InstructPix2Pix [2], where users specify the classifier-free guidance weights in order to explore the desired amount of change for a given edit. In our method, we inherit this parameter. Below we list these chosen values for the sequences shown in the paper:

1. *Fig. 4 Tree*:  $s_I=1.5, s_T \in [7.5, 12.5]$
2. *Fig. 8 Tents*:  $s_I=1.5, s_T=10.0$
3. *Fig. 4 Person*:  $s_I \in [1.5, 1.75], s_T \in [6.5, 7.5]$
4. *Fig. 1*:  $s_I=1.5, s_T \in [6.5, 7.5]$
5. *Fig. 4 Bear*:  $s_I = 1.5, s_T = 6.5$
6. *Fig. 7*:  $s_I \in [1.3, 1.5], s_T \in [6.5, 8.5]$

Our process of optimizing for an edited NeRF uses a diffusion model as guidance, which can produce a collection of temporally varying images (i.e., varying over the course of optimization). As a result, the optimization process does not have a single convergence point, as standard NeRF optimization does, where all images are sufficiently well explained by the reconstructed model. Therefore, the edited NeRF also varies in type and strength of edit over the course of optimization, and one must select an iteration at which to terminate optimization and visualize the edited scene. In practice, the optimal choice for training length is a subjective decision — a user may prefer more subtle or more extreme edits that are best found at different stages of training. The results in the paper are shown after a varying number of iterations, 3000 – 4000 for smaller scenes (e.g., Fig. 1 and Fig. 4 bear, which both have under 100 images each), and 7000 – 8000 for larger scenes (the remaining scenes, with over 200-300 images each).

For InstructPix2Pix, we use the public implementation included in the Diffusers library [45].

## B. Limitations

As mentioned in the main paper, our method inherits many of the limitations of InstructPix2Pix. This includes (1) the inability to perform large spatial manipulations, (2) the occasional inability to perform binding between objects referred to in the instruction and the corresponding scene objects, and (3) adding or removing large objects. Furthermore, as in DreamFusion, our method uses a diffusion



Figure 10: **Effects of autoencoder**: From left to right: (1) the input captured image, (2) the original reconstructed NeRF, rendered from the same viewpoint, (3), the original image 1, encoded and decoded by Stable Diffusion’s autoencoder, (4) a NeRF reconstruction of the autoencoded images, rendered from the same viewpoint. Below, we show a zoomed-in patch of the eye. We note that while the input sequence can be reasonably reconstructed without much loss of detail, and the autoencoded images are similarly sharp as the input captured sequence, there are slight perturbations in local textures that end up not being 3D-consistent, resulting in blurry renders in the final reconstruction.

model on a single view at a time, and thus may suffer from similar artifacts, such as double faces on added objects. Finally, it’s worth noting that the edit instructions provided to InstructPix2Pix are sometimes more relevant to certain views than others. For example, if the instruction is to “turn the man into a bear”, not all views may prominently feature the man, and therefore, certain views may consistently produce less of an edit or no edit at all. Our framework can easily incorporate improvements made in the InstructPix2Pix and its follow up works.

We additionally note that the edited NeRF scenes often contain slightly blurrier textures when compared to the originally reconstructed NeRF. Some initial experimentation suggests that this may be a result of the Stable Diffusion autoencoder, which often does not produce an entirely faithful copy of the original image, even in unedited regions. The autoencoder, while producing visually comparable results, often creates images with locally similar but non-identical textures that are not globally 3D-consistent. To validate this hypothesis, we show an experiment in Figure 10, where we simply autoencode the input images using Stable Diffusion (i.e., the same autoencoder used by InstructPix2Pix) and continue training the NeRF. As a result, the scene becomes gradually blurrier.

Furthermore, over much longer optimization runs (i.e., when the process of rendering, editing, and propagating the edited images back into the NeRF has been repeated many times), we note a decrease in visual quality. We also largely

attribute this to the effects of the autoencoder.

## C. Metrics

In the quantitative evaluation, we report two metrics, a *CLIP Directional Score* [6], and a metric we name *CLIP Direction Consistency Score*. The CLIP Directional score measures how much the change in text captions agrees with the change in the images, and the CLIP Consistency score measures the cosine similarity of the CLIP embeddings of each pair of adjacent frames in a render novel camera path.

More formally, for the directional score, we encode a pair of images (the original and edited NeRFs, rendered at a given viewpoint), as well as a pair of text prompts that describe the original and edited scenes, e.g., “*a photograph of a man*” and “*a photograph of a Tolkien Elf*”. Using these, we compute the directional score described in InstructPix2Pix [2] and StyleGAN-NADA [6].

For the temporal consistency loss, we encode a pair of consecutive rendered frames from a novel trajectory, using both the original NeRF and the edited NeRF. In all, we are left with four CLIP embeddings,  $C(o_i)$ ,  $C(o_{i+1})$ ,  $C(e_i)$ ,  $C(e_{i+1})$ , corresponding to original NeRF renderings  $o_i, o_{i+1}$  and edited NeRF renderings  $e_i, e_{i+1}$  for consecutive novel views  $i$  and  $i + 1$ . We define the consistency loss as:

$$(C(e_i) - C(o_i)) \cdot (C(e_{i+1}) - C(e_i)) \quad (2)$$

This effectively amounts to measuring the change in the CLIP-space edit direction from frame to frame.