

CLIP-NeRF: Text-and-Image Driven Manipulation of Neural Radiance Fields

Can Wang¹, Menglei Chai², Mingming He³, Dongdong Chen⁴, Jing Liao^{1,†}

¹City University of Hong Kong ²Snap Inc. ³USC Institute for Creative Technologies ⁴Microsoft Cloud AI

cwang355-c@my.cityu.edu.hk, cmlatsim@gmail.com, hmm.lillian@gmail.com

cddlyf@gmail.com, jingliao@cityu.edu.hk

Abstract

We present CLIP-NeRF, a multi-modal 3D object manipulation method for neural radiance fields (NeRF). By leveraging the joint language-image embedding space of the recent Contrastive Language-Image Pre-Training (CLIP) model, we propose a unified framework that allows manipulating NeRF in a user-friendly way, using either a short text prompt or an exemplar image. Specifically, to combine the novel view synthesis capability of NeRF and the controllable manipulation ability of latent representations from generative models, we introduce a disentangled conditional NeRF architecture that allows individual control over both shape and appearance. This is achieved by performing the shape conditioning via applying a learned deformation field to the positional encoding and deferring color conditioning to the volumetric rendering stage. To bridge this disentangled latent representation to the CLIP embedding, we design two code mappers that take a CLIP embedding as input and update the latent codes to reflect the targeted editing. The mappers are trained with a CLIP-based matching loss to ensure the manipulation accuracy. Furthermore, we propose an inverse optimization method that accurately projects an input image to the latent codes for manipulation to enable editing on real images. We evaluate our approach by extensive experiments on a variety of text prompts and exemplar images and also provide an intuitive interface for interactive editing.

1. Introduction

With the explosive growth of 3D assets, the demand for manipulating 3D content to achieve versatile re-creation is rising rapidly. While most existing 3D editing methods operate on explicit 3D representations [7, 15, 45], the recent advances of implicit volumetric representations in capturing

and rendering dedicated 3D structures [9, 14, 16, 23, 27, 34] have motivated the research to benefit the manipulation from such representations. Among these works, neural radiance fields (NeRF) [23] utilize a volume rendering technique to render neural implicit representations for high-quality novel view synthesis, providing an ideal representation for 3D content.

Editing NeRF (*e.g.*, deforming the shape or changing the appearance color), however, is an extremely challenging task. First, since NeRF is an implicit function optimized per scene, we cannot directly edit the shape using the intuitively tools for the explicit representations [35, 40–42]. Second, unlike image manipulation where the single-view information is enough to guide the editing [20, 43, 44], the multi-view dependency of NeRF makes the manipulation way more difficult to control without the multi-view information. More recent works propose conditional NeRF [36], which trains NeRF on one category of shapes and enables manipulation via latent space interpolations utilizing the pre-trained models. Based on the conditional NeRF, Edit-NeRF [21] takes the first step to edit the shape and color of NeRF given user scribbles. However, due to its limited capacity in shape manipulation, only adding or removing local parts of the object is allowed. In addition to achieving more compelling and complicated manipulation, we seek to edit NeRF in more intuitive ways, such as using a text prompt or a single reference image.

In this paper, we explore how to individually manipulate the shape and the appearance of NeRF based on a text prompt or a reference image in a unified framework. Our framework is built on a novel disentangled conditional NeRF architecture, which is controlled by the latent space disentangled into a shape code and an appearance code. The shape code guides the learning of a deformation field to warp the volume to a new geometry, while the appearance code allows controlling the emitted color of volumetric rendering. Based on our disentangled NeRF model, we take advantage of the recently proposed Contrastive Language-Image Pre-training (CLIP) model [33] to learn two code mappers, which map CLIP features to the latent space to

[†]Jing Liao is the corresponding author. Our project page is [h
//cassiepython.github.io/clipnerf/](https://cassiepython.github.io/clipnerf/)

manipulate the shape or appearance code. Specifically, given a text prompt or an exemplar image as our condition, we extract the features using the pre-trained CLIP model, feed the features into the code mappers, and yield local displacements in the latent space to edit the shape and appearance codes to reflect the edit. We design the CLIP-based loss to enforce the CLIP space consistency between the input constraint and the output renderings, thus supporting high-resolution NeRF manipulation. Additionally, we propose an optimization-based method for editing a real image by inversely optimizing its shape and appearance codes.

To sum up, we make the following contributions:

- We present the first text-and-image-driven manipulation method for NeRF, using a unified framework to provide users with flexible control over 3D content using either a text prompt or an exemplar image.
- We design a disentangled conditional NeRF architecture by introducing a shape code to deform the volumetric field and an appearance code to control the emitted colors.
- Our feedforward code mappers enable the fast inference for editing different objects in the same category compared to the optimization-based editing method [21].
- We propose an inversion method to infer the shape and appearance codes from a real image, allowing editing the shape and appearance of the existing data.

2. Related Work

NeRF and NeRF Editing. The past few years have witnessed tremendous progress in the implicit representation of 3D models with neural networks [9, 14, 16, 23, 27, 34]. Among them, NeRF [23] is a representative one, which encodes a continuous volume representation of shape and view-dependent appearance in the weights of an MLP network. NeRF has been gaining more and more popularity because of its strong capability in capturing high-resolution geometry and rendering photo-realistically novel views. The success of NeRF has also inspired many follow-up works that extend the NeRF to dynamic scenes [8, 29, 32, 39], relighting [3, 38], generative models [4, 13, 24, 36], etc. Furthermore, DietNeRF [12] designs a CLIP semantic consistency loss to improve few-shot NeRF and presents impressive results, and GRAF [36] first adopts shape and appearance codes to conditionally synthesize NeRF, which inspires our adversarial training.

Despite the above success, a 3D model with NeRF representation is very unintuitive and difficult to edit since it is represented by millions of network parameters. To address this problem, the pioneering work EditNeRF [21] designs a

conditional NeRF, where the 3D object encoded by NeRF is conditioned on a shape code and an appearance code. By optimizing the adjustment to these two latent codes, user edits on shape and appearance color can be achieved. However, this method has limited capacity in shape manipulation as it only supports adding or removing local parts of the object. Also, the editing process of EditNeRF [21] is slow because of its iterative optimization nature. Compared to EditNeRF [21], our method is different in three aspects. First, our method gives more freedom in shape manipulation and supports global deformation. Second, by learning two feed-forward networks mapping user edits to the latent codes, our method allows fast inference for the interactive editing. Moreover, different from the user scribbles used in EditNeRF [21], we introduce two intuitive ways to NeRF editing: using either a short text prompt or an exemplar image, which are more friendly to novice users.

CLIP-Driven Image Generation and Manipulation. An important building block of our work is CLIP [33] which connects texts and images by bringing them closer in a shared latent space, under a contrastive learning manner. Powered by the CLIP model, some text-driven image generation and manipulation methods are proposed. Perez [31] combines CLIP and StyleGAN [17, 18] to synthesize images by optimizing the latent code of a pre-trained StyleGAN according to a textual condition defined in the CLIP space. Instead of generating images from scratch, StyleCLIP [30] introduces a text-based interface for StyleGAN to allow manipulations of real images with text prompts. Besides applying CLIP to GAN models, DiffusionCLIP [19] combines a diffusion model [37] with CLIP to conduct a text-driven image manipulation. It achieves a comparable performance to that of GAN-based image manipulation methods, with the advantage of great mode coverage and training stability. However, all these methods only explore the text-guidance ability of CLIP, whereas our method unifies both text-and-image driven manipulations in a single model by fully exploiting the power of CLIP. Further, these methods are limited to image manipulation and fail to encourage multi-view consistency due to the lack of 3D information. In contrast, our model combines NeRF with CLIP, thus allowing editing 3D models in a view consistent way.

3. Method

In this section, we start with the general formulation of conditional NeRF (§ 3.1) as a 3D generative model conditioned by shape and appearance codes. We then present our disentangled conditional NeRF model (§ 3.2), which is able to individually control the shape and appearance manipulation. Next, we introduce our framework on leveraging the multi-modal power of CLIP for driving NeRF manipulation (§ 3.2.2) using both text prompts or image exemplars, and the inversion strategy (§ 3.4). Finally, we propose an inversion

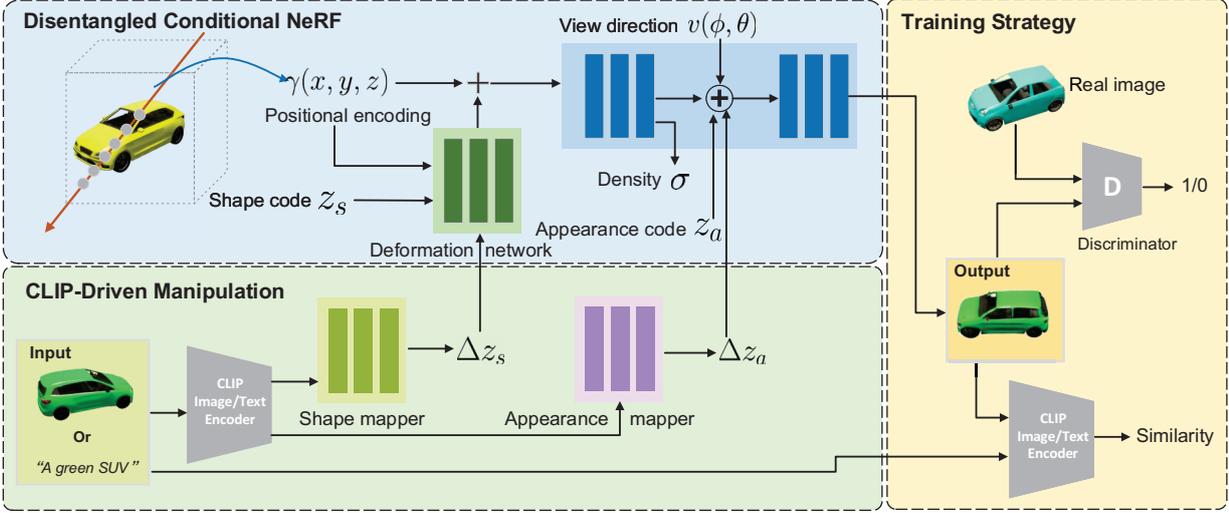


Figure 1. **The framework of the proposed method.** Our model first learns a disentangled conditional NeRF which takes positional encoding, view direction, shape code, and appearance code as input and outputs rendered image, while the shape code aims to deform the volume field via a deformation network. This disentangled conditional NeRF is trained in an adversarial manner. Then given a reference image or a text prompt, the CLIP image or text encoder extracts the corresponding feature embedding for the shape and appearance mappers to learn a local step in the latent space for shape and appearance manipulation, respectively. These two mappers are trained using a CLIP similarity loss with our pre-trained disentangled conditional NeRF.

method (§ 3.5) to allow editing a real image by a novel latent optimization approach on shape and appearance codes.

3.1. Conditional NeRF

Built upon the original per-scene NeRF, conditional NeRF servers as a generative model for a particular object category, conditioned on the latent vectors that dedicatedly control shape and appearance. Specifically, conditional NeRF is represented as a continuous volumetric function \mathcal{F}_θ that maps a 5D coordinate (a spatial position $\mathbf{x}(x, y, z)$ and a view direction $\mathbf{v}(\phi, \theta)$), together with a shape code \mathbf{z}_s and an appearance code \mathbf{z}_a , to a volumetric density σ and a view-dependent radiance $\mathbf{c}(r, g, b)$, parametrized by a multi-layer perceptron (MLP). A trivial formulation $\mathcal{F}'_\theta(\cdot)$ of conditional NeRF can be:

$$\mathcal{F}'_\theta(\mathbf{x}, \mathbf{v}, \mathbf{z}_s, \mathbf{z}_a) : (\Gamma(\mathbf{x}) \oplus \mathbf{z}_s, \Gamma(\mathbf{v}) \oplus \mathbf{z}_a) \rightarrow (\mathbf{c}, \sigma), \quad (1)$$

where \oplus is the concatenation operator.

Here $\Gamma(\mathbf{p}) = \{\gamma(p) \mid p \in \mathbf{p}\}$ is the sinusoidal positional encoding that separately projects each coordinate p of vector \mathbf{p} to a high dimensional space. Each output dimension of the encoding function $\gamma(\cdot) : \mathbb{R} \rightarrow \mathbb{R}^{2m}$ is defined as:

$$\gamma(p)_k = \begin{cases} \sin(2^k \pi p), & \text{if } k \text{ is even,} \\ \cos(2^k \pi p), & \text{if } k \text{ is odd,} \end{cases} \quad (2)$$

where $k \in \{0, \dots, 2m-1\}$ and m is a hyper-parameter that controls the total number of frequency bands.

3.2. Disentangled Conditional NeRF

The aforementioned conditional NeRF does introduce conditional generation capability to the NeRF architecture. However, this trivial formulation \mathcal{F}'_θ (Eq. 1) suffers from mutual intervention between shape and appearance conditions, e.g., manipulating the shape code could also cause color changes. In observation of this issue, we propose our disentangled conditional NeRF architecture to achieve individual control over both shape and appearance by properly disentangling the conditioning mechanism.

Conditional Shape Deformation. Rather than directly concatenating the latent shape code to the encoded position feature, we propose to formulate the shape conditioning through explicit volumetric deformation to the input position. This conditional shape deformation not only improves the robustness of the manipulation and preserves the original shape details as much as possible by regularizing the output shape to be a smooth deformation of the base shape, but more importantly also completely isolates the shape condition from affecting the appearance.

To this end, we design a shape deformation network $\mathcal{T} : (\mathbf{x}, \mathbf{z}_s) \rightarrow \Delta \mathbf{x}$, which projects a position \mathbf{x} and the input \mathbf{z}_s to displacement vectors $\Delta \mathbf{x} \in \mathbb{R}^{3 \times 2m}$ corresponding to each band of the positional encoding $\Gamma(\mathbf{x})$. Thus, the deformed positional encoding $\Gamma^*(\mathbf{p}, \mathbf{z}_s) = \{\gamma^*(p, \Delta p) \mid p \in \mathbf{p}, \Delta p \in \mathcal{T}(\mathbf{p}, \mathbf{z}_s)\}$ is defined as:

$$\gamma^*(p, \Delta p)_k = \gamma(p)_k + \tanh(\Delta p_k), \quad (3)$$

where the scalar p and the vector $\Delta p \in \mathbb{R}^{2m}$ belong to the k -th axis from \mathbf{p} and $\Delta \mathbf{p}$. The hyperbolic tangent function

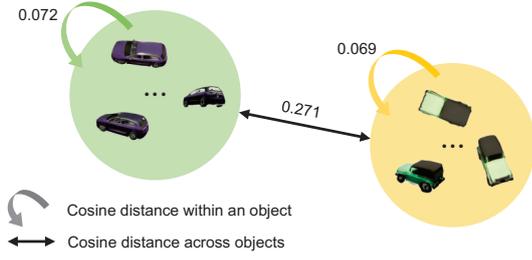


Figure 2. **Multi-View Consistency Evaluation of CLIP.** We randomly select two cars and measure their pairwise CLIP cosine distances of 1) different cars in a same view and 2) a same car in different random views (we sampled 144 views from the upper hemisphere, as a combination of 12 ϕ and 12 θ poses). Though the camera poses vary dramatically, different views for a same object have higher similarity (small distance). But different objects have lower similarity (large distance) even in an identical view.

$\tanh(\cdot)$ is used to constrain the displacements in the range of $[-1, 1]$, which helps avoid poor local minimums caused by large motions and increase the training robustness.

Deferred Appearance Conditioning. In NeRF, the density is predicted first as a function of position and the radiance is then predicted from both position and view direction. Similar to Graf [36] and EditNeRF [21], we also defer the appearance conditioning to concatenate the appearance code with the view direction as the input to the radiance prediction network, which allows manipulating the appearance without touching the shape information, i.e., density.

Overall, as illustrated in Fig. 1, our disentangled conditional NeRF $\mathcal{F}_\theta(\cdot)$ is defined as:

$$\mathcal{F}_\theta(\mathbf{x}, \mathbf{v}, \mathbf{z}_s, \mathbf{z}_a) : (\Gamma^*(\mathbf{x}, \mathbf{z}_s), \Gamma(\mathbf{v}) \oplus \mathbf{z}_a) \rightarrow (\mathbf{c}, \sigma). \quad (4)$$

And for the simplicity of notation, we use $\mathcal{F}_\theta(\mathbf{v}, \mathbf{z}_s, \mathbf{z}_a) = \{\mathcal{F}_\theta(\mathbf{x}, \mathbf{v}, \mathbf{z}_s, \mathbf{z}_a) \mid \mathbf{x} \in \mathbf{R}\}$ to denote the rendering of the whole image with viewport \mathbf{R} .

3.3. CLIP-Driven Manipulation

With our disentangled conditional NeRF (Eq. 4) as a generator, we now introduce how we integrate the CLIP model into the pipeline to achieve text-driven manipulation on both shape and appearance.

To avoid optimizing both shape and appearance codes for each target sample, which tends to be versatile and time-consuming, we take a feed-forward approach to directly update the condition codes from the input text prompt. Specifically, given an input text prompt of \mathbf{t} and the initial shape/appearance code of $\mathbf{z}'_s/\mathbf{z}'_a$, we train a shape mapper \mathcal{M}_s and an appearance mapper \mathcal{M}_a to update the codes as:

$$\begin{aligned} \mathbf{z}_s &= \mathcal{M}_s(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}'_s, \\ \mathbf{z}_a &= \mathcal{M}_a(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}'_a, \end{aligned}$$

3838

where $\hat{\mathcal{E}}_t(\cdot)$ is the pre-trained CLIP text encoder that projects the text to the CLIP embedded feature space and both mappers map this CLIP embedding to displacement vectors that update the original shape and appearance codes.

In addition, given that CLIP includes an image encoder and a text encoder mapping to a joint embedding space, we define a cross-modal CLIP distance function $D_{\text{CLIP}}(\cdot, \cdot)$ to measure the embedding similarity between the input text and a rendered image patch:

$$D_{\text{CLIP}}(\mathbf{I}, \mathbf{t}) = 1 - \langle \hat{\mathcal{E}}_i(\mathbf{I}), \hat{\mathcal{E}}_t(\mathbf{t}) \rangle, \quad (6)$$

where $\hat{\mathcal{E}}_i(\cdot)$ and $\hat{\mathcal{E}}_t(\cdot)$ are the pre-trained CLIP image and text encoders, \mathbf{I} and \mathbf{t} are the input image patch and text, and $\langle \cdot, \cdot \rangle$ is the cosine similarity operator.

Without loss of generality, here we assume that the manipulation control comes from a text prompt \mathbf{t} . However, our distance can also be extended to measure similarity between two images or two text prompts. Thus, our framework naturally supports editing with an image exemplar by trivially replacing the text prompt with this exemplar in aforementioned equations.

Discussion. To perform NeRF manipulation with image-level CLIP model, a natural question is whether the CLIP feature is stable across different viewpoints and whether it can distinguish object differences. To evaluate this, we randomly select two objects (e.g., an SUV and a jeep) and measure the pairwise CLIP-space cosine distances between 1) different views of a same object, and 2) different objects in a same view. As shown in Fig. 2, we find the distance is more sensitive to small object difference than large view variations. This suggests that a pre-trained CLIP model has the ability to support view-consistency representations for 3D-aware applications. A similar observation is found by DietNeRF [12] and applied in 3D reconstruction.

3.4. Training Strategy

Our pipeline is trained in two stages: we first train the disentangled conditional NeRF including the conditional NeRF generator and the deformation network; then we fix the weights of the generator and train the CLIP manipulation parts including both the shape and appearance mappers.

Disentangled Conditional NeRF. Our conditional NeRF generator \mathcal{F}_θ is trained together with the deformation network using a non-saturating GAN objective [22] with the discriminator \mathcal{D} , where $f(x) = -\log(1 + \exp(-x))$ and λ_r is the regularization weight. Assuming that real images \mathbf{I} form the training data distribution of d , we randomly sample the shape code \mathbf{z}_s , the appearance code \mathbf{z}_a , and the camera pose from \mathcal{Z}_s , \mathcal{Z}_a , and \mathcal{Z}_v , respectively, where \mathcal{Z}_s and \mathcal{Z}_a are normal distribution, and \mathcal{Z}_v is the upper hemisphere

of the camera coordinate system.

$$\mathcal{L}_{\text{GAN}} = \mathbb{E}_{\mathbf{z}_s \sim \mathcal{Z}_s, \mathbf{z}_a \sim \mathcal{Z}_a, \mathbf{v} \sim \mathcal{Z}_v} [f(\mathcal{D}(\mathcal{F}_\theta(\mathbf{v}, \mathbf{z}_s, \mathbf{z}_a)))] + \mathbb{E}_{\mathbf{I} \sim \mathcal{d}} [f(-\mathcal{D}(\mathbf{I}) + \lambda_r \|\nabla \mathcal{D}(\mathbf{I})\|^2)]. \quad (7)$$

CLIP Manipulation Mappers. We use pre-trained NeRF generator \mathcal{F}_θ , CLIP encoders $\{\hat{\mathcal{E}}_t, \hat{\mathcal{E}}_i\}$, and the discriminator \mathcal{D} to train the CLIP shape mapper \mathcal{M}_s and appearance mapper \mathcal{M}_a . All network weights, except the mappers, are fixed, denoted as $\{\hat{\cdot}\}$. Similar to the first stage, we randomly sample the shape code \mathbf{z}_s , the appearance code \mathbf{z}_a , and the camera pose \mathbf{v} from their respective distributions. In addition, we sample the text prompt \mathbf{t} from a pre-defined text library \mathbf{T} . By using our CLIP distance D_{CLIP} (Eq. 6) with weight λ_c , we train the mappers with the following losses:

$$\mathcal{L}_{\text{shape}} = f(\hat{\mathcal{D}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \mathcal{M}_s(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}_s, \mathbf{z}_a))) + \lambda_c D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \mathcal{M}_s(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}_s, \mathbf{z}_a), \mathbf{t}), \quad (8)$$

$$\mathcal{L}_{\text{appear}} = f(\hat{\mathcal{D}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \mathbf{z}_s, \mathcal{M}_a(\hat{\mathcal{E}}_i(\mathbf{t})) + \mathbf{z}_a))) + \lambda_c D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \mathbf{z}_s, \mathcal{M}_a(\hat{\mathcal{E}}_i(\mathbf{t})) + \mathbf{z}_a), \mathbf{t}). \quad (9)$$

3.5. Inverse Manipulation

The manipulation pipeline we have introduced so far works on an initial sample with known conditions including the shape and appearance codes. To apply the manipulation to an input image \mathbf{I}_r belonging to the same training category, the key is to first optimize all generation conditions to inversely project the image to the generation manifold, similar to the latent image manipulation methods [1, 2, 10, 26]. Following the EM algorithm [5], we design an iterative method to alternatively optimize the shape code \mathbf{z}_s , the appearance code \mathbf{z}_a , and the camera \mathbf{v} . To be specific, during each iteration, we first optimize \mathbf{v} while keeping \mathbf{z}_s and \mathbf{z}_a fixed using the following loss:

$$\mathcal{L}_v = \|\hat{\mathcal{F}}_\theta(\mathbf{v}, \hat{\mathbf{z}}_s, \hat{\mathbf{z}}_a) - \mathbf{I}_r\|_2 + \lambda_v D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \hat{\mathbf{z}}_s, \hat{\mathbf{z}}_a), \mathbf{I}_r). \quad (10)$$

We then update the shape code by minimizing:

$$\mathcal{L}_s = \|\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \mathbf{z}_s + \lambda_n \mathbf{z}_n, \hat{\mathbf{z}}_a) - \mathbf{I}_r\|_2 + \lambda_s D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \mathbf{z}_s + \lambda_n \mathbf{z}_n, \hat{\mathbf{z}}_a), \mathbf{I}_r), \quad (11)$$

where \mathbf{z}_a and \mathbf{v} are fixed, \mathbf{z}_n is a random standard Gaussian noise vector sampled in each step to improve the optimization robustness, and λ_n linearly decays from 1 to 0 through the whole optimization iterations.

The appearance code is updated in a similar manner:

$$\mathcal{L}_a = \|\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \hat{\mathbf{z}}_s, \mathbf{z}_a + \lambda_n \mathbf{z}_n) - \mathbf{I}_r\|_2 + \lambda_a D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \hat{\mathbf{z}}_s, \mathbf{z}_a + \lambda_n \mathbf{z}_n), \mathbf{I}_r), \quad (12)$$

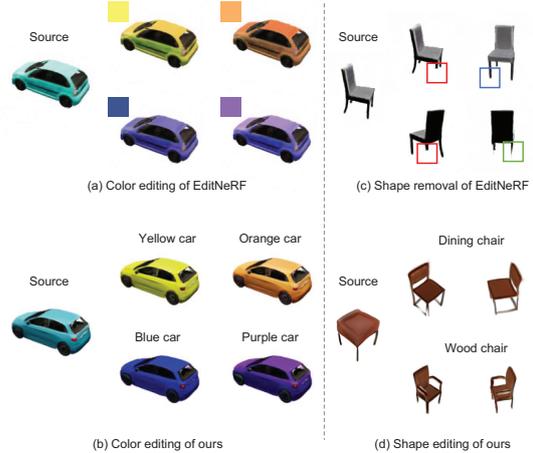


Figure 3. Compared to EditNeRF.

	Chairs		Cars	
	Shape	Appearance	Shape	appearance
EditNeRF	30.0	15.9	33.2	16.8
Ours	0.58	0.51	2.12	1.98

Table 1. **Compared to EditNeRF [21] on editing time averaged on 20 images.** We only include the inference/optimization time(s) and single-view rendered time(s) for chairs (128×128 pixels) and cars (256×256 pixels).

4. Experiments

Datasets. We evaluate our method on two public datasets: *Photoshapes* [28, 36] with 150K chairs rendered at 128×128 following the rendering protocol of [25] and *Carla* with 10K cars rendered at 256×256 using the Driving simulator [6, 36]. Each object is rendered in a random view without providing any camera pose parameters.

Implementation details. Our conditional NeRF is an 8-layer MLP with each layer containing 256 hidden units, and the input dimension is 64. Following the default architecture for NeRF [23], we also use ReLU activations. The deformation network is a 4-layer MLP with ReLU activations and 256 hidden units per layer. It takes a 128-dimensional shape code \mathbf{z}_s as input, $\mathbf{z}_s \in \mathbb{R}^{128}$. We also represent the appearance code \mathbf{z}_a using 128 dimensions, $\mathbf{z}_a \in \mathbb{R}^{128}$. Both shape and appearance mappers are 2-layer MLPs with ReLU activations. The channel sizes of each mapper are 128, 256, and 128, respectively. The implementation of the discriminator follows PatchGAN [11]. We use the Adam optimizer and an initial learning rate of 10^{-4} to train the network. The learning rate is decayed by 0.5 every 50K steps. In the inversion, we also use the Adam optimizer with the learning rate starting from 10^{-3} and decreasing by 0.75 every 100 steps. Besides, $\lambda_r = 0.5$, $\lambda_v = 0.1$, and $\lambda_s - \lambda_a = 0.2$. All the models are trained on an NVIDIA 3839 GPU platform.



Figure 4. **Ablation Study for Disentanglement.** We show text-and-exemplar driven shape editing results of our method and the baseline method without using our disentangled technique. When editing the shape, the latter can change the appearance, while ours keeps the appearance unchanged.

4.1. Compared to EditNeRF

We compare with pioneering work in NeRF editing, EditNeRF [21] on the editing of shape and appearance color of both datasets in Fig. 3. For the Photoshapes dataset, EditNeRF is trained using 600 instances with 40 views per instance while ours uses only one view. For the Carla dataset, EditNeRF uses 10K cars with a single view per instance, same as ours. Besides, camera pose parameters are required during training of EditNeRF but unknown for us.

We first compare the capability and performance between EditNeRF and our method. For the color editing (Fig. 3-(a)), EditNeRF requires the user to select a target color and draw coarse scribbles on a local region. With the foreground and background masks created by the coarse scribbles, EditNeRF performs the appearance editing by optimizing the appearance code and a conditional NeRF to achieve the target color. We observe that unnatural color effects appear on the edited results of EditNeRF (*e.g.* discontinuity on car doors), and the generated color is not completely faithful to the target color. In contrast, we allow the user to change the color more simply by providing a text prompt and our method produces more natural editing results (Fig. 3-(b)). For the shaping editing (Fig. 3-(c)), EditNeRF can only support local shape editing, such as shape part removal. Given the user’s editing scribbles which for example indicate to remove a leg of a chair (in the red rectangles), EditNeRF optimizes a few layers in the network to fit the shape in the input view but it cannot ensure successful propagation to unseen views (in the blue rectangle) and keep the structure of other parts intact (in the green rectangle). Compared to it, our method supports a large degree of shape deformation and generalizes well to unseen views (Fig. 3-(d)). Besides, EditNe



Figure 5. **Ablation study on our inversion method and comparison with EditNeRF.**

an optimized-based method, takes a large amount of time for the optimization, while our feedforward code mappers achieve much faster inference of the target shape and appearance (Table 1).

To quantitatively evaluate how good the image quality is preserved after editing, we calculate the FID scores of 2K testing images before and after editing. Due to training with 40 views per instance, EditNeRF shows better reconstruction before editing on the chair dataset, but its editing notably degrades the image quality while our method ensures comparable quality before and after manipulation. On the car dataset, the performance of EditNeRF significantly drops because of only one view per instance used in training. Trained under the same setting, our model improves reconstruction quality by a large margin and well preserves the quality during editing. Since EditNeRF requires user scribbles for shape editing and is difficult to generate a large set of results with random conditions, we exclude it in the comparison on shape editing while our method performs equally well regardless of editing shape or color.

We also compare with EditNeRF in the inversion results (Fig. 5). EditNeRF infers the shape and appearance codes by fine-tuning the condition NeRF using the standard NeRF photometric loss. Our optimized-based inversion method outperforms by benefiting from CLIP’s ability to provide multi-view consistency representations (more discussions in Section 3.3 and ablation in 4.2).

4.2. Ablation Study

We evaluate our model w/ and w/o the disentangled design (Section 3.2). In Fig. 4, the model trained without the conditional shape deformation network (*i.e.* w/o disen.) frequently introduces color changes when performing shape editing. In contrast, our disentangled conditional NeRF achieves individual shape control because the conditional deformation network is able to isolate the shape con-

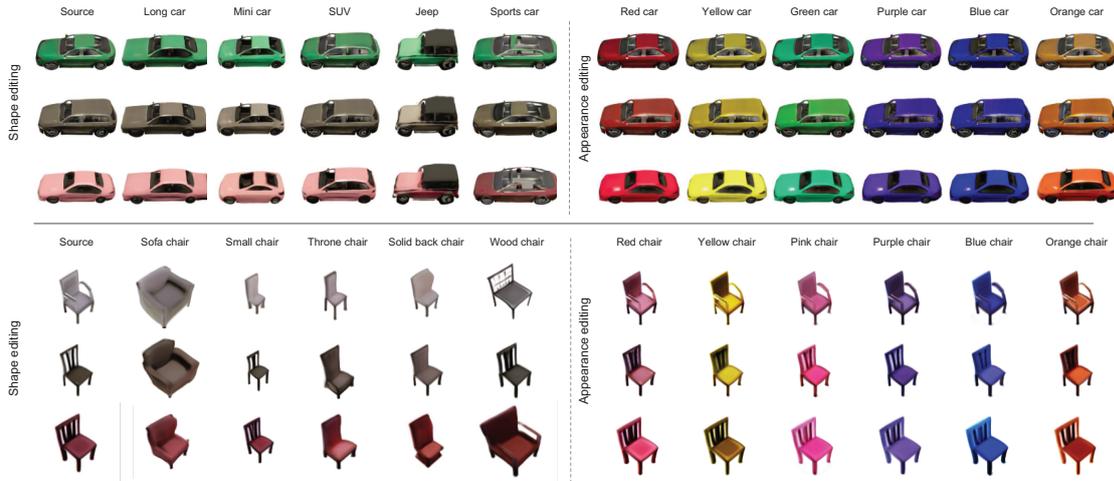


Figure 6. Text-Driven Editing Results.

	Chairs			Cars		
	Before	After	Diff.	Before	After	Diff.
(a) EditNeRF	36.8	40.2	3.4	102.8	118.7	15.9
w/o disen.	52.5	54.3	1.8	69.2	69.9	0.7
Ours	47.8	49.0	1.2	66.7	67.2	0.5
(b) w/o disen.	52.5	53.2	0.7	69.2	71.1	1.9
Ours	47.8	48.4	0.6	66.7	67.8	1.1

Table 2. Fréchet inception distance (FID) for evaluating the image quality of reconstructed views before and after editing on: (a) color and (b) shape (lower value means better). We use 2K images with various views drawn randomly from the latent space to calculate the FIDs for reconstructed images, and then perform various edits on these images to recalculate FIDs of edited results. As EditNeRF requires user scribbles for shape editing, we exclude it from the comparison on shape manipulation. w/o disen. is a variant of our model without the shape deformation network used for disentangling control of shape and appearance.

dition from the appearance control and deform the base volume field to generate new objects without affecting the appearance. Also, since the deformation network implicitly enforces regularization of the generated shape, the resulting quality is further improved as shown in Table 4.2.

We conduct another ablation study in Fig. 5 to evaluate our inversion optimization method. The baseline method (w/o CLIP) only computes the standard NeRF photometric loss between the output and a single image. Its result quality is limited due to the difficulty in inferring a complete 3D NeRF model from a single view. As discussed in Section 3.3, CLIP has the capability to produce robust pose-invariant features. Therefore, our inversion method introduces a CLIP constraint during optimization and achieves better inversion results thanks to the CLIP prior.

4.3. CLIP-Driven Manipulation

Our method supports editing of object shape or appearance using text. When manipulating the shape, we keep

appearance code unchanged and the same applies to appearance editing. Fig. 6 demonstrates diverse editing results. Note that when the car with a light color is deformed to a sports car, its color may become darker. But it is not a failure case, as the colors of all sports cars in the Carla dataset are inherently intenser. Besides, we find that our method naturally preserves the shading when changing the appearance color. When editing the chair shape, if the user’s input text is highly relevant to the source shape—for example, the source chair is a wood chair, and the user also wants a ‘wood chair’—the result will be slightly different from the source. During the color editing, our method guarantees that the shape is completely preserved. Our method also supports exemplar-based manipulation by providing a real target image instead of a text prompt. We present various exemplar-guided shape and appearance editing results in Fig. 7. Our method achieves semantic-precise and individual control of shape and appearance referring to the exemplar image.

4.4. Real Image Manipulation

To evaluate the generalization ability of our model in processing a single real image that does not exist in our training set, we experiment with the real image by inverting it to a shape code and an appearance code and then applying them to edit. We show the inverted and edited results in Fig. 8. We observe that inverting the chair is much more challenging than inverting a car due to the chair’s delicate structures, such as the wheels of the office chair. However, even the office chair is not perfectly reconstructed, the editing ability of our method is not affected. Our method still ensures accurate editing in shape and appearance.

4.5. User Study

We conduct a user study to evaluate the perceptual quality and accuracy of the editing results. We include 20 ques-

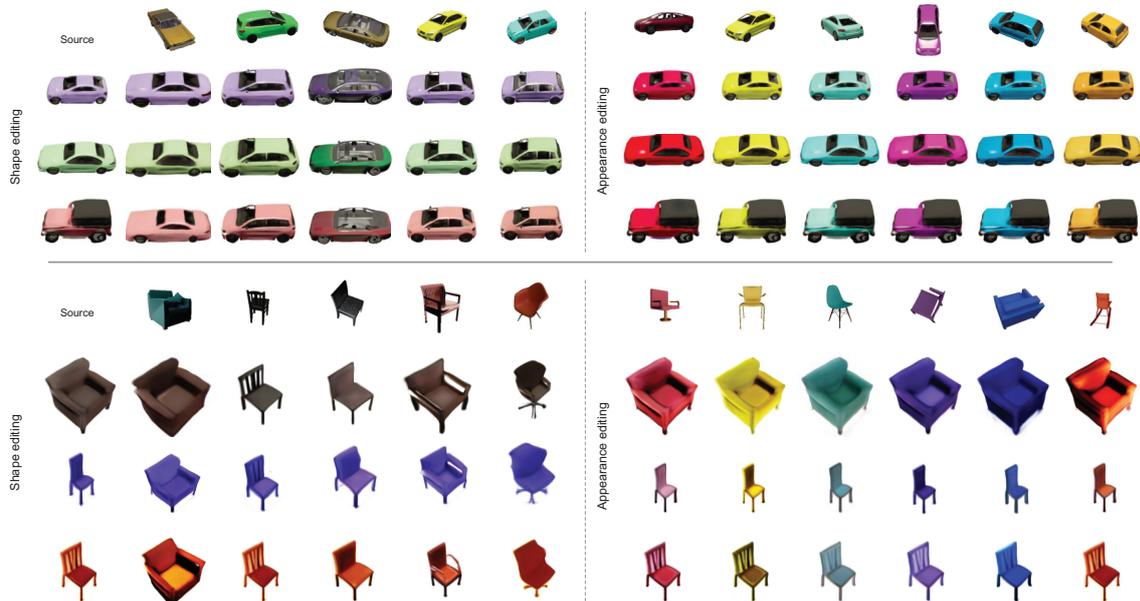


Figure 7. Exemplar-Driven Editing Results.

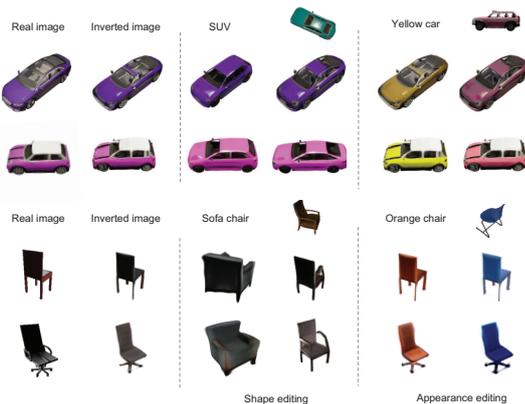


Figure 8. Editing Results on Real Images.

Chairs			Cars		
Text	Exemplar	Avg.	Text	Exemplar	Avg.
0.821	0.877	0.849	0.814	0.859	0.837

Table 3. **User Study Results.** We report correct matching rate by counting whether an editing result is matched to the right text/exemplar guidance by users.

tions in the study, each question with 5 results of cars or chairs generated by 5 randomly selected text prompts or 5 randomly selected exemplars. We randomly shuffle the results and give users unlimited time to match each result with the correct text or image. We collect answers from 23 participants and report the matching accuracy rate in Table 3. Our method, in more than 80% cases, succeeds in editing objects exactly corresponding to the description given by the text or the exemplar.

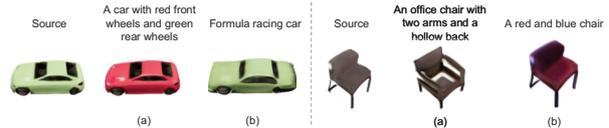


Figure 9. **Limitations.** Our method cannot handle fine-grained edits (a) and out-of-domain edits (b).

5. Conclusion

We present the first text-and-image driven manipulation method for NeRF by designing a unified framework to provide users with flexible control over 3D content using either a text prompt or an exemplar image. We design a disentangled conditional NeRF architecture that allows disentangling shape and appearance while editing an object, and two feedforward code mappers enable fast inference for editing different objects. Further, we proposed an inversion method to infer the shape and appearance codes from a real image, allowing editing the existing data.

Limitations. We evaluate our approach by extensive experiments on various text prompts and exemplar images and provide an intuitive editing interface for interactive editing. However, our method cannot handle fine-grained and out-of-domain shape and appearance edits as shown in Fig. 9, due to the limited expressive ability of the latent space and the pre-trained CLIP. This may be alleviated by adding more various training data.

Acknowledgement This work was supported by the Hong Kong Innovation and Technology Commission (ITC) under the Innovation and Technology Fund (Project No. MHP/109/19).

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4432–4441, 2019. 5
- [2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embedded images? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8296–8305, 2020. 5
- [3] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. Nerd: Neural reflectance decomposition from image collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12684–12694, 2021. 2
- [4] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5799–5809, 2021. 2
- [5] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977. 5
- [6] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. In *Conference on robot learning*, pages 1–16. PMLR, 2017. 5
- [7] Ohad Fried, Ayush Tewari, Michael Zollhöfer, Adam Finkelstein, Eli Shechtman, Dan B Goldman, Kyle Genova, Zeyu Jin, Christian Theobalt, and Maneesh Agrawala. Text-based editing of talking-head video. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 1
- [8] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2021. 2
- [9] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4857–4866, 2020. 1, 2
- [10] Jinjin Gu, Yujun Shen, and Bolei Zhou. Image processing using multi-code gan prior. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3012–3021, 2020. 5
- [11] 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. 5
- [12] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, 2021. 2, 4
- [13] Wonbong Jang and Lourdes Agapito. Codenerf: Disentangled neural radiance fields for object categories. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12949–12958, 2021. 2
- [14] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, Thomas Funkhouser, et al. Local implicit grid representations for 3d scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6001–6010, 2020. 1, 2
- [15] Tao Ju, Qian-Yi Zhou, and Shi-Min Hu. Editing the topology of 3d models by sketching. *ACM Transactions on Graphics (TOG)*, 26(3):42–es, 2007. 1
- [16] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 364–375, 2017. 1, 2
- [17] 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
- [18] 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
- [19] Gwanghyun Kim and Jong Chul Ye. Diffusionclip: Text-guided image manipulation using diffusion models. *arXiv preprint arXiv:2110.02711*, 2021. 2
- [20] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip HS Torr. Manigan: Text-guided image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7880–7889, 2020. 1
- [21] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. Editing conditional radiance fields. *arXiv preprint arXiv:2105.06466*, 2021. 1, 2, 4, 5, 6
- [22] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *International conference on machine learning*, pages 3481–3490. PMLR, 2018. 4
- [23] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. 1, 2, 5
- [24] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. 2
- [25] Michael Oechsle, Michael Niemeyer, Christian Reiser, Lars Mescheder, Thilo Strauss, and Andreas Geiger. Learning implicit surface light fields. In *2020 International Conference on 3D Vision (3DV)*, pages 452–462. IEEE, 2020. 5
- [26] Xingang Pan, Xiaoang Zhan, Bo Dai, Dahua Lin, Chen Change Loy, and Ping Luo. Exploiting deep generative prior for versatile image restoration and manipulation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 5

- [27] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 1, 2
- [28] Keunhong Park, Konstantinos Rematas, Ali Farhadi, and Steven M Seitz. Photoshape: photorealistic materials for large-scale shape collections. *ACM Transactions on Graphics (TOG)*, 37(6):1–12, 2018. 5
- [29] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Deformable neural radiance fields. *arXiv preprint arXiv:2011.12948*, 2020. 2
- [30] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2085–2094, 2021. 2
- [31] Victor Perez. Generating images from prompts using CLIP and StyleGAN. <https://towardsdatascience.com/generating-images-from-prompts-using-clip-and-stylegan-1f9ed495ddda>, 2021. 2
- [32] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. 2
- [33] 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. *arXiv preprint arXiv:2103.00020*, 2021. 1, 2
- [34] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *European Conference on Computer Vision*, pages 623–640. Springer, 2020. 1, 2
- [35] Thorsten-Walther Schmidt, Fabio Pellacini, Derek Nowrouzezahrai, Wojciech Jarosz, and Carsten Dachsbacher. State of the art in artistic editing of appearance, lighting and material. In *Computer Graphics Forum*, volume 35, pages 216–233. Wiley Online Library, 2016. 1
- [36] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *arXiv preprint arXiv:2007.02442*, 2020. 1, 2, 4, 5
- [37] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Proceedings of the 33rd Annual Conference on Neural Information Processing Systems*, 2019. 2
- [38] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7495–7504, 2021. 2
- [39] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhofer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12959–12970, 2021. 2
- [40] Mikaela Angelina Uy, Jingwei Huang, Minhyuk Sung, Tolga Birdal, and Leonidas Guibas. Deformation-aware 3d model embedding and retrieval. In *European Conference on Computer Vision*, pages 397–413. Springer, 2020. 1
- [41] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Hang Yu, Wei Liu, Xiangyang Xue, and Yu-Gang Jiang. Pixel2mesh: 3d mesh model generation via image guided deformation. *IEEE transactions on pattern analysis and machine intelligence*, 2020. 1
- [42] Weiyue Wang, Duygu Ceylan, Radomir Mech, and Ulrich Neumann. 3dn: 3d deformation network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1038–1046, 2019. 1
- [43] Weihao Xia, Yujiu Yang, Jing-Hao Xue, and Baoyuan Wu. Tedigan: Text-guided diverse face image generation and manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2256–2265, 2021. 1
- [44] Xiaogang Xu, Ying-Cong Chen, Xin Tao, and Jiaya Jia. Text-guided human image manipulation via image-text shared space. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 1
- [45] Matthias Zwicker, Mark Pauly, Oliver Knoll, and Markus Gross. Pointshop 3d: An interactive system for point-based surface editing. *ACM Transactions on Graphics (TOG)*, 21(3):322–329, 2002. 1