# PaletteNeRF: Palette-based Color Editing for NeRFs

Qiling Wu
Tsinghua University
Qinghua Yuan, Haidian, Beijing
wql19@mails.tsinghua.edu.cn

Jianchao Tan
Kuaishou Technology
No.6 Shangdi West Road, Haidian, Beijing
tanjianchaoustc@gmail.com

Kun Xu
Tsinghua University
Qinghua Yuan, Haidian, Beijing
xukun@tsinghua.edu.cn

## Abstract

**Neural Radiance Field (NeRF) is a powerful tool to faithfully generate novel views for scenes with only sparse captured images. Despite its strong capability for representing 3D scenes and their appearance, its editing ability is very limited. In this paper, we propose a simple but effective extension of vanilla NeRF, named *PaletteNeRF*, to enable efficient color editing on NeRF-represented scenes. Motivated by recent palette-based image decomposition works, we approximate each pixel color as a sum of palette colors modulated by additive weights. Instead of predicting pixel colors as in vanilla NeRFs, our method predicts additive weights. The underlying NeRF backbone could also be replaced with more recent NeRF models such as KiloNeRF to achieve real-time editing. Experimental results demonstrate that our method achieves efficient, view-consistent, and artifact-free color editing on a wide range of NeRF-represented scenes.**

*Keywords: NeRF, NeRF editing, palette, recoloring*

## 1. Introduction

Neural Radiance Field (NeRF) [22] is a powerful tool for image-based modeling and rendering. With only a sparse set of captured images, it can faithfully generate rendering results from novel views. The core of NeRF is a neural volumetric representation of the scene using a multi-layer perceptron network. Due to its high effectiveness, it has attracted a wide range of attention from the community, with a bunch of follow-up works and applications since its introduction in 2020. However, due to the black-box nature of neural representations, all information, including geometries, materials, and light transports, are tightly baked into NeRFs, which are hard to be interpreted and edited.
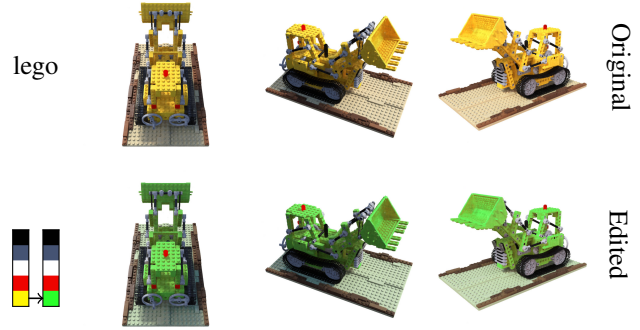


Figure 1: Color editing results of scene *lego* represented by Neural Radiance Field (NeRF). The origin and edited palettes are shown at the bottom-left corner. Users can modify the extracted palette to achieve intuitive, view-consistent, artifact-free editing of NeRFs.

Recently, some methods have been proposed to enable color editing of NeRFs [20, 42]. Given a few user-provided color scribbles, EditingNeRF [20] propagates the user edits to the whole data to achieve color editing and shape modification. However, this work is too demanding for datasets, which require many instances from the same category for training, limiting its practical usage. CLIP-NeRF [42] uses embeddings of CLIP [32] to edit NeRFs. They finetune layers that influence color while freezing layers that influence density, to match the embedding of NeRF's output to that of the text editing prompt. However, it sometimes modifies undesired areas, leaving some artifacts in the results. PosterNeRF [40] gives an efficient way to extract palette from radiance fields, then utilizes posterization method [5] to achieve real-time color editing. Despite the real-time performance, the editing results have artifacts like color banding and leaking as side effects of posterization.

In this paper, we propose PaletteNeRF, an intuitive,

1

view-consistent and artifact-free palette-based color editing method for NeRFs. Inspired by palette-based image editing works, we approximate pixel colors as the sum of palette colors modulated by additive weights. Instead of predicting pixel colors as in vanilla NeRFs, we predict the high dimensional additive weights in our PaletteNeRF. The NeRF-represented scenes can be recolored by adjusting the palette colors without retraining or modifying the PaletteNeRF, as shown in Figure 1. The underlying NeRF backbone could also be replaced by more recent NeRF models such as Kilo-NeRF [33] to achieve real-time editing.

## 2. Related work

### 2.1. Neural Radiance Field

Neural Radiance Field (NeRF) [22] utilizes radiance field to model a 3D scene implicitly. Specifically, they use MLPs to infer volumetric density and radiance for certain points and view directions of a scene and follow the paradigm of volumetric rendering to compute the image pixel colors. Given a sparse set of captured images, NeRF can generate high-quality results for novel views. However, this framework needs to train a separate MLP for every scene. All information about the scene, including geometries, materials, and light transports, is baked into the neural representation. Hence, the vanilla NeRF does not allow changes in scene geometries, colors, and lighting. Various follow-ups of NeRFs have been proposed to address this limitation, i.e., to deal with deformable scenes [31, 41, 28, 29], dynamic lighting [21], and scene composition [25, 47, 46, 44, 14, 27].

To improve NeRF's rendering speed, NSVF [19] and KiloNeRF [33] use empty space skipping and early ray termination. In addition, KiloNeRF divides the scene into small grids, and uses a small and efficient MLP in each grid to achieve further speedup. FastNeRF [12] and PlenOctrees [45] utilize function factorization and cache the MLP results to speedup rendering. Some other methods aim at efficient training, e.g., Instant-NGP [23] use hierarchical hash encoding to replace the costly MLP.

As for editing, EditingNeRF [20] takes a set of objects with similar shapes and colors, such as cars from Carla dataset [8], then associate each instance with shape code and color code, which are fed into Conditional Radiance Fields (CRFs). The users provide a few color scribbles to indicate color changes and shape addition/removal of an object, which are propagated to specific regions of that object. However, it requires a set of shapes in the same class with different geometries and colors. CLIP-NeRF [42] takes texts or exemplar images as edit prompts, which are fed into the multi-modal language model, i.e. CLIP, to obtain an embedding to serve as an offset for shape and color codes. The CRF process the modified codes to output edited re-

sults. PosterNeRF [40] extracts palette efficiently from radiance fields. They sample RGB color points and use volumetric visibility to remove outliers, after which, a palette is extracted from the remaining points. Each pixel is approximated as a linear blending of at most 2 palette colors. Color editing is performed by adjusting the palette color. However, undesired artifacts can be easily perceived in the editing results. In contrast, our work only requires a single scene to enable color editing. Our method provides a relatively simple user interface and achieves intuitive, artifact-free, and view-consistent color editing. Readers could refer to more details in the surveys on NeRFs [6, 39].

### 2.2. Palette based editing

A palette is a concise representation of the color distributions of an image or video, which can be used to efficiently edit the image or video. Several works have studied the human perception of the palette. These works construct various data fitting models according to human preference, and then regress a perceptual palette from input images [26, 18, 3, 10]. Other works [4, 24, 48, 1] adopt clustering methods (e.g., k-means) over pixel colors, and use the cluster centers as palette colors. Another type of works [38, 37, 43, 7, 13, 16] utilize geometric methods to generate palettes for images. Specifically, Tan et al. [38, 37] compute an RGB space convex hull that includes all colors of an image, followed by an iterative simplification of the convex hull until the vertex number is reduced to a predefined number or the reconstruction error reaches a predefined threshold. The vertices of the simplified convex hull are considered palette colors. Recently, Du et al. [9] further extend this geometry-based method to time-lapse videos, generating time-varying palettes. With the help of a palette, an image can be decomposed into multiple compositing layers, each of which contains spatially-variant per-pixel opacities or per-pixel mixing weights. Several methods [35, 38] generate ordered translucent layers by assuming an alpha blending color composition mode. Other methods [2, 17, 37, 43, 34] generate order-independent layers by assuming additive color composition mode.

## 3. Background

### 3.1. Neural Radiance Field (NeRF)

Given a sparse set of captured images on a scene, NeRF [22] can faithfully synthesize novel views for the scene. The core is a neural volumetric representation modeled with a multi-layer perceptron (MLP). The MLP $f$ predicts the RGB color and density observed at a 3D position $\mathbf{x}$ from a view direction $\mathbf{d}$:

$$(\mathbf{c}, \sigma) = f(\gamma(\mathbf{x}), \gamma(\mathbf{d})), \tag{1}$$

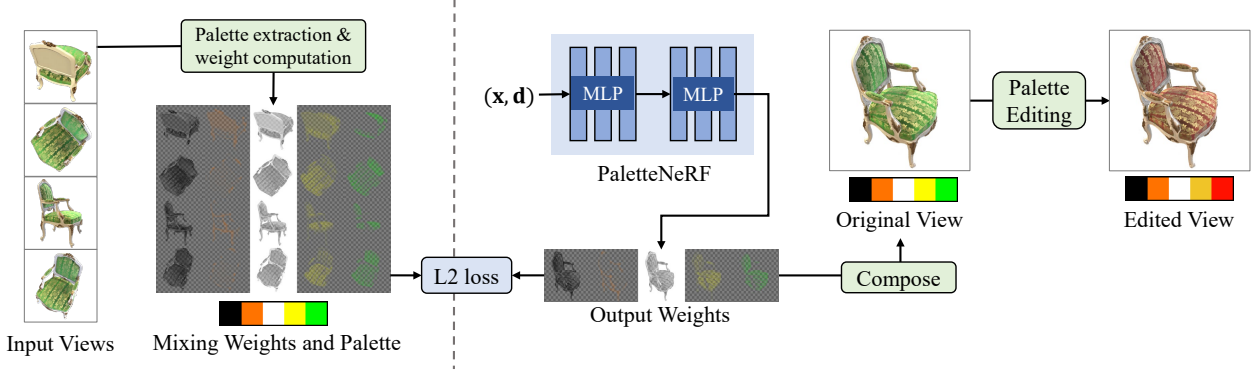where $\gamma(\cdot)$ denotes a positional encoding function.

Figure 2: The pipeline of PaletteNeRF. We first preprocess input views to obtain the palette and mixing weights. The weights are used as NeRF's fitting target. After training NeRF, we obtain output weights for a novel view, and modify the palette color to edit that view.

To render an image from a novel view, NeRF follows conventional ray marching techniques for volumetric rendering. For each image pixel, we first cast a ray $\mathbf{r}$ from the viewpoint through the pixel to the scene. After that, several 3D points are sampled along the ray, and each sampled point is fed into the MLP to obtain its color and density. Finally, the colors of all sampled points are blended to obtain the rendered pixel color.

A sparse set of captured images together with their camera parameters is sufficient for training a NeRF. A separate NeRF is required to be trained for each separate scene. During training, at each iteration, a batch of rays is randomly sampled from all (or a subset of) pixels. L2 differences between rendered pixel colors and ground truth pixel colors at given views are used as the loss function to be minimized.

### 3.2. Palette based image decomposition

Given an RGB image $\boldsymbol{I}$, palette based image decomposition techniques [38, 37, 43, 9] extract a small set of colors named *palette* from an image, such that the color $\mathbf{c}$ of each pixel $\mathbf{p}$ could be approximately represented as a linear combination of palette colors:

$$\mathbf{c} \approx \sum_{1 \leq i \leq K} w_i^{\mathbf{P}} \cdot \mathbf{v}_i, \qquad (2)$$

where $K$ denotes palette size (i.e., number of palette colors) which usually varies from 4 to 10, depending on the color complexity of the input image. $\mathbf{v}_i$ denotes the $i$-th palette color, and $w_i^{\mathbf{P}}$ denotes the additive weight of pixel $\mathbf{p}$ with respect to palette color $\mathbf{v}_i$. The pixel-wise additive weights are required to satisfy two properties. First, they are non-negative: $w_i^{\mathbf{P}} \geq 0$ ($1 \leq i \leq K$); second, they sum to one: $\sum_{i=1}^{K} w_i^{\mathbf{P}} = 1$.

By denoting the palette as $\boldsymbol{V} = [\mathbf{v}_1, \ldots, \mathbf{v}_K]$ and the additive weight vector at each pixel $\mathbf{p}$ as $\mathbf{w}^{\mathbf{P}} =$

$[w_1^{\mathbf{P}}, \ldots, w_K^{\mathbf{P}}]^T$, respectively, we could also rewrite the above formula in the vector-matrix form as:

$$\mathbf{c} \approx \boldsymbol{V} \cdot \mathbf{w}^{\mathbf{P}}. \qquad (3)$$

By denoting the per-pixel additive weights ($\mathbf{w}^{\mathbf{P}}$) as a weight image $\boldsymbol{W}$, i.e., having the same resolution as the input image $\boldsymbol{I}$ while the channel size is changed from 3 to $K$, we could approximate pixel colors of the input image by:

$$\boldsymbol{I} \approx \boldsymbol{V} \cdot \boldsymbol{W}. \qquad (4)$$

There are various options of additive weights, including as-sparse-as-possible weights [38], RGBXY barycentric weights [37], and mean value coordinates (MVC) [43, 15, 11]. Since MVC weights could be efficiently computed and have been demonstrated to be smooth, sparse and effective [43, 9], we use MVC weights in our paper.

Once the palette and the additive weight images are extracted, by simply adjusting the colors in the palette, the images can be instantly recolored through weighted linear interpolations from the palette colors (Eq. 2). This offers a more convenient editing interface than stroke-based methods, since those methods additionally user interaction steps, i.e., drawing strokes on the image.

## 4. Our Method

In this section, we introduce PaletteNeRF, a concise modification to NeRF that makes scenes editable. We will first introduce our pipeline (Sec. 4.1), then discuss the details of our network structure and training strategy (Sec. 4.2).

### 4.1. Overview and Pipeline

Our goal is to provide an intuitive, efficient, and artifact-free color editing tool for NeRFs. Motivated by existing
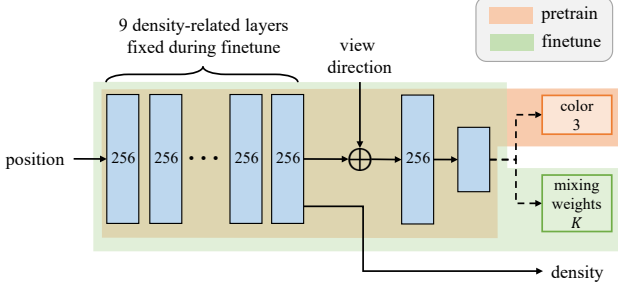
Figure 3: The network structure of our PaletteNeRF. The number inside each block denotes the dimension of the layer. The top branch of the network predicts 3D colors which are used in the pretrain stage, while the bottom branch predicts $K$-channel additive mixing weights which are used in the finetune stage.

palette-based image decomposition works [38, 37, 43, 9], we find that if we approximate additive weights rather than pixel colors, the colors of the scene represented by NeRFs could be naturally edited.

Specifically, instead of using NeRFs to model pixel colors which are essentially 3D (RGB) signals, we model higher-dimensional signals, e.g., the pixel-wise additive weights. Thanks to the high capacities of MLPs, the $K$-dimensional pixel-wise additive weights could also be well approximated. After that, we could intuitively edit the colors of the NeRF-represented scene by adjusting the palette colors. We refer to our modified NeRF as *PaletteNeRF*.

Our overall pipeline is shown in Figure 2. The steps for utilizing PaletteNeRF include:

- **Data Preparation.** Given a set of sparse captured RGB images of a scene, we convert them to the same number of $K$-channel additive weight images and a shared palette $V$ with $K$ colors. This is done by stitching all images into a big image, and then directly applying the method in [43]. The pixel-wise additive weights are obtained using mean value coordinates.

- **Modified MLP.** We feed the $K$-channel additive weight images to PaletteNeRF. Different from the vanilla NeRF which predicts 3-channel RGB colors, our PaletteNeRF predicts a $K$-channel vector $\mathbf{w}$ indicating additive weights, i.e., its formulation is changed from $f(\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma)$ to $f(\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{w}, \sigma)$.

- **Novel View Rendering.** To obtain the additive weight image $W_\mathbf{o}$ from a novel view, we apply the same process as done in vanilla NeRF, i.e., through ray marching and blending of $K$-channel values. The RGB image from this novel view can be simply reconstructed through Eq. 2 as:

$$I_\mathbf{o} = V \cdot W_\mathbf{o}. \tag{5}$$

- **Color Editing.** Then, we can freely adjust the palette colors $V$ to achieve recoloring of the scene. During the color editing process, our PaletteNeRF, which predicts additive weights, is not modified at all. This indicates that our editing process is decoupled with the Palette-NeRF structure. Such decoupling has two advantages. First, the used NeRF backbone could be replaced with any newer, faster variants of NeRFs. In our experiments, we also demonstrate our method with a Kilo-NeRF backbone [33], which allows color editing in real-time. Second, the consistency between different views is also naturally preserved, leading to artifact-free color editing with high fidelity.

### 4.2. Network Structure and Training

As shown in Figure 3, our network structure is almost the same as the vanilla NeRF, except for the last layer. We use a $128 \times K$ fully connected (FC) layer, instead of a $128 \times 3$ FC layer, as the last layer to output a $K$-channel vector.

For training a PaletteNeRF, we need to change the loss function from computing L2 differences between rendered and GT RGB images to computing L2 differences between rendered (predicted) and the GT $K$-channel additive weight images. A straightforward training process would be directly using the same process as done for training a vanilla NeRF, only with the above change on the loss function. While such a straightforward training strategy could work, we find that its reconstruction quality is relatively low.

Based on the fact that our PaletteNeRF and the vanilla NeRF share a majority part of network structures, we propose a two-stage *pretrain + finetune* training strategy. Specifically, we first pretrain a vanilla NeRF using RGB input images. After that, we fix the parameters of the density-related layers, i.e. the parameters in the first 9 layers are kept unchanged, and finetune the parameters of the last 2 color-related layers by minimizing L2 differences between rendered and the GT $K$-channel additive weight images. Experiments in Sec. 6.1 show that such a two-stage training strategy leads to better reconstruction accuracy.

## 5. Method Extensions

Our method can be easily extended in several ways. In the following, first, we show that the backbone of Palette-NeRF can be replaced by KiloNeRF [33] to achieve real-time recoloring (Sec. 5.1). Second, we extend the palette and mixing weights of PaletteNeRF to capture indirect lighting in synthetic scenes (Sec. 5.2). Finally, we would like to show that applications in palette-based editing can be directly incorporated into our framework. Specifically, we demonstrate an application of color harmonization with PaletteNeRF (Sec. 5.3).
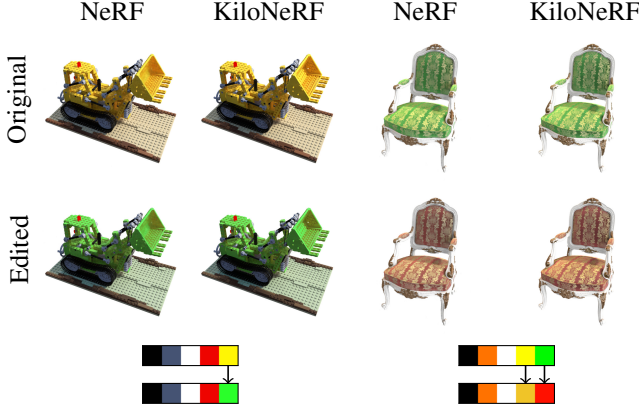
Figure 4: Editing results of using vanilla NeRF and Kilo-NeRF as backbones, respectively. Both methods can fit the multi dimensional weights well, and achieve editing results with visually indistinguishable quality.

## 5.1. Real-time Color Editing

As mentioned earlier, our editing process is decoupled with the underlying NeRF backbone. The vanilla NeRF backbone could be replaced by any newer, faster variants of NeRFs as long as they are still volumetric representations. To demonstrate this ability, we also implement our Palette-NeRF with a KiloNeRF backbone [33], which allows color editing in real time. Similar to NeRF, we modify the last FC layer of the small MLP used in KiloNeRF from $32 \times 3$ to $32 \times K$. This only slows down the rendering time for each frame by around 4% for $K = 5$ and still achieves real-time framerates. The editing results of the KiloNeRF backbone are shown in Figure 4. The results with KiloNeRF backbones are visually indistinguishable from those with vanilla NeRF backbones.

## 5.2. Second-order weights for synthetic scenes

NeRF scenes could be classified into captured scenes and synthetic scenes. Captured scenes refer to those scenes where only captured images from specific views are available. While the images of synthetic scenes are generated using renderers with given 3D geometries, materials, and lighting information.

For synthetic NeRF scenes, we could further extend our method to provide a more semantically meaningful palette and better controls on recoloring, i.e., supporting higher-order effects (indirect lighting) besides additive mixing. Let's imagine a synthetic 3D scene with static geometry and static lighting, but with some editable materials, whose diffuse colors $\mathbf{v}_i$ are allowed to adjust. Considering indirect illuminations up to 2 bounces, the rendered color $\mathbf{c}$ of a pixel
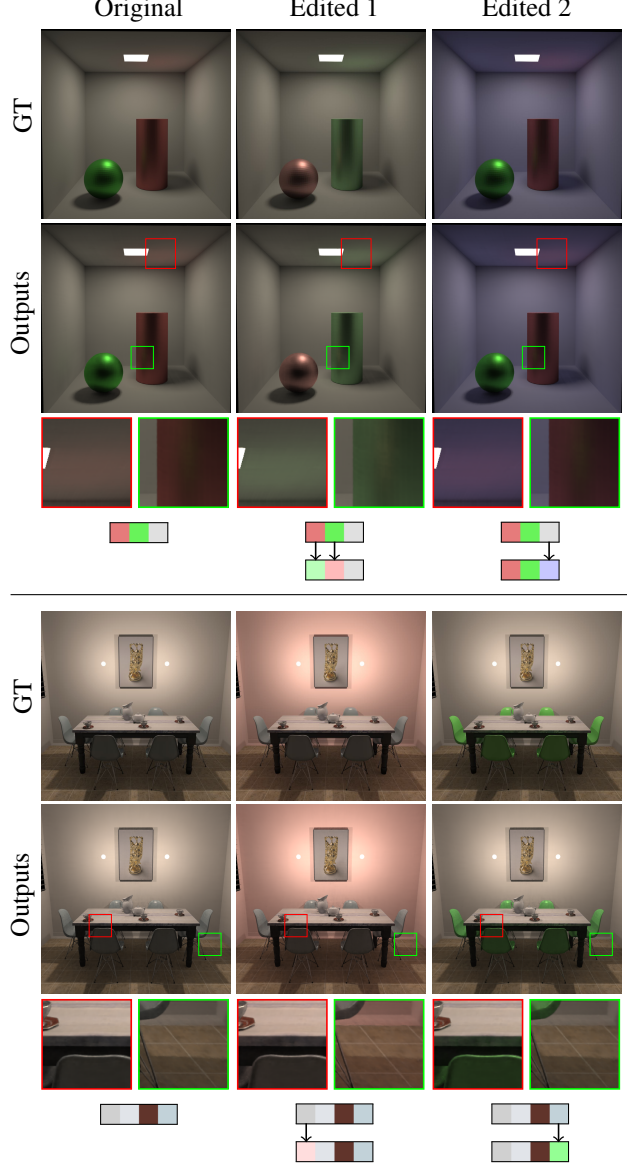


Figure 5: Rendered scenes (*Cornell box, breakfast*) and editing results. For each scene, the 4 rows are: ground truth, PaletteNeRF output, amplified areas, color of tunable materials. Colored rectangles amplify area which have interreflection effects. Palette colors are gamma-corrected.

$\mathbf{p}$ could be computed as follows:

$$\mathbf{c} = w_0^{\mathbf{P}} + \sum_{1 \leq i \leq K} w_i^{\mathbf{P}} \cdot \mathbf{v}_i + \sum_{1 \leq i \leq j \leq K} w_{i,j}^{\mathbf{P}} \cdot (\mathbf{v}_i \odot \mathbf{v}_j), \ (6)$$

where $\odot$ denotes channel-wise multiplication, $w_0^{\mathbf{P}}$ denotes contributions from light paths not intersecting with any editable materials. $w_i^{\mathbf{P}}$ and $w_{i,j}^{\mathbf{P}}$ denote contributions from light paths intersecting one and two times of editable materials, respectively. All these weights are essentially light

transports and can be directly computed using a path tracing renderer such as PBRT [30].

The equation can be written in the vectorized form:

$$\mathbf{c} = \boldsymbol{V} \cdot \mathbf{w^P}, \tag{7}$$

where:

$$\begin{aligned}
\boldsymbol{V} &= [\mathbf{1}, \mathbf{v}_1, \dots, \mathbf{v}_K, \mathbf{v}_{1,1}, \dots, \mathbf{v}_{K,K}], \\
\mathbf{w^P} &= [w_0^{\mathbf{P}}, w_1^{\mathbf{P}}, \dots, w_K^{\mathbf{P}}, w_{1,1}^{\mathbf{P}}, \dots, w_{K,K}^{\mathbf{P}}]^T.
\end{aligned} \tag{8}$$

The above modification could be directly supported by our PaletteNeRF, by only modifying the data preparation step at the beginning and the color editing step at the end, as described in Sec. 4.1. Notice that a palette size of $K$ would produce a second-order weight vector $\boldsymbol{V}$ of length $(K + 1)(K + 2)/2$.

In the data preparation step, we no longer use existing palette-based image decomposition methods to extract palette and additive weights. Instead, we ask users to manually specify several materials he or she wants to adjust and use the diffuse colors of specified materials as palette colors. Then, we compute the second-order weights in Eq. 8 directly using PBRT [30]. In the color editing step, we use Eq. 6 instead of Eq. 2 to compute the reconstructed colors.

We have tested two synthetic scenes: *Cornell box* and *breakfast*. The editing results are shown in Figure 5. Notice the visually plausible indirect illumination effects after recoloring, e.g., color bleeding effects in the ceiling and reflection of the sphere on the cylinder surface in scene *Cornell box*, and color bleeding effects including pink tint of the floor and green tint of the desk in scene *breakfast*.

### 5.3. Color Harmonization

Our method could be directly combined with the palette-based image harmonization method [36] to achieve color harmonization for NeRFs.

Specifically, we use the 7 harmonic color templates $T_m, (m = 1 \dots 7)$ designed for palette colors in [36]. Each harmonic template has at least 1 parameter $\alpha$, which describes the angle of rotation in the hue of the axis. The harmonization is performed in HCL color space (Hue, Chroma, Luminance). We utilize the palette $\boldsymbol{V}$ extracted in the data preparation step of PaletteNeRF. For template $T_m(\alpha)$, we find the closest axis to each color in $\boldsymbol{V}$, measured by the distance of hue in HCL color space. The distance is weighted and summed to compute a template fitting cost, which is then minimized to find the best $\alpha$ for each template.

Once the template is chosen, each palette color in $\boldsymbol{V}$ is projected onto the nearest axis of the template, forming the harmonized palette $\boldsymbol{V}'$. Then, as described in Sec. 4.1, $\boldsymbol{V}'$ is used to recolor, i.e. harmonize the whole scene. Figure 6 shows harmonized results of scene *orchids*.
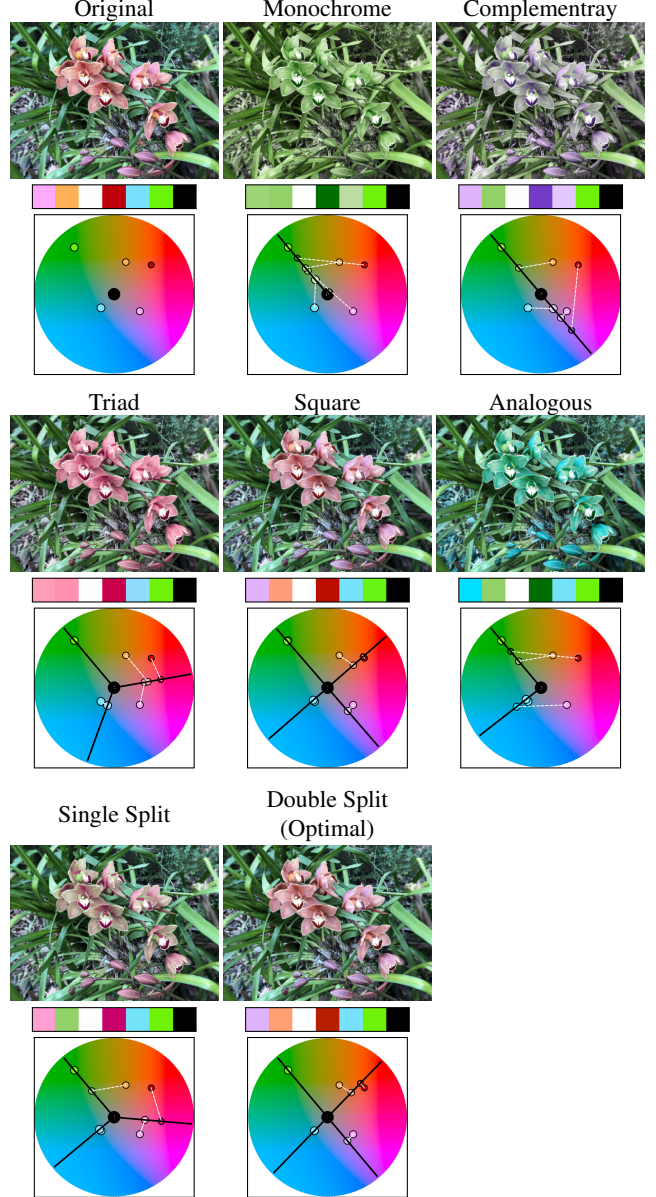


Figure 6: Comparison of the same scene *orchids* harmonized by 7 different templates.

## 6. Experiments

We conduct our experiments on a PC with an NVIDIA RTX 3080Ti GPU, a Ryzen 5900X CPU, and 64GB of RAM. For each scene, we train PaletteNeRF for 200k iterations, which take around 5 hours.

### 6.1. Evaluation

**Training strategy.** As mentioned in Sec. 4.2, we have proposed a two-stage (pretrain + finetune) training strategy. To validate the effectiveness of the two-stage strategy, we

| training | Additive Weights | | | Second Order Weights | | |
|---|---|---|---|---|---|---|
| method | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| traditional | 25.205 | 0.832 | 0.129 | 8.579 | 0.084 | 0.363 |
| ours | 26.711 | 0.833 | 0.095 | 40.107 | 0.976 | 0.024 |

Table 1: Quantitative comparison between the traditional strategy and our two-stage (pretrain + finetune) training strategy. "Additive Weights" columns report average scores of 6 captured scenes (*lego, chair, flower, orchids, vasedeck, valley*). "Second Order Weights" columns report average scores of 2 synthetic scenes (*Cornell box, breakfast*).

| | vasedeck | | | lego | | |
|---|---|---|---|---|---|---|
| dimension | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| 128 | 24.262 | 0.693 | 0.172 | 30.297 | 0.941 | 0.031 |
| 256 | 24.730 | 0.722 | 0.154 | 30.525 | 0.944 | 0.030 |
| 512 | 24.732 | 0.720 | 0.156 | 30.664 | 0.945 | 0.029 |
| NeRF(128) | 24.713 | 0.715 | 0.163 | 30.604 | 0.946 | 0.030 |

Table 2: Ablation studies on the dimension of the last layer of PaletteNeRF (shown in the first 3 rows) evaluated on scenes *vasedeck* and *lego*. The last row shows scores of the vanilla NeRF with 128 dimension for comparison.

compare it against the traditional one. For a fair comparison, we set the number of iterations to be the same for both training strategies. Specifically, the traditional strategy uses 200k iterations, while our two-stage strategy set iteration steps as 180k (pretrain) + 20k (finetune) for captured scenes, and 140k (pretrain) + 60k (finetune) for synthetic scenes. We have tested 8 scenes, including 6 captured scenes and 2 synthetic scenes, and report the average reconstruction scores for both strategies in Table 1.

From the results, we could find that our two-stage training strategy improves the reconstruction quality by a large gap, especially on the synthetic scenes with second-order weights. This is possibly due to the relatively large number of channels making it harder to be directly trained, while a pretrain step helps provide a good initialization.

**Dimension of the last layer.** We also evaluate different choices of the input dimension of the last layer. Table 2 shows the reconstruction scores on two scenes when the dimension of the last layer is set to 128, 256, and 512, respectively. Increasing the dimension of the last layer could slightly increase reconstruction quality, at the cost of longer rendering time. However, to make a good trade-off between quality and speed, we set it as 128 in our experiments.
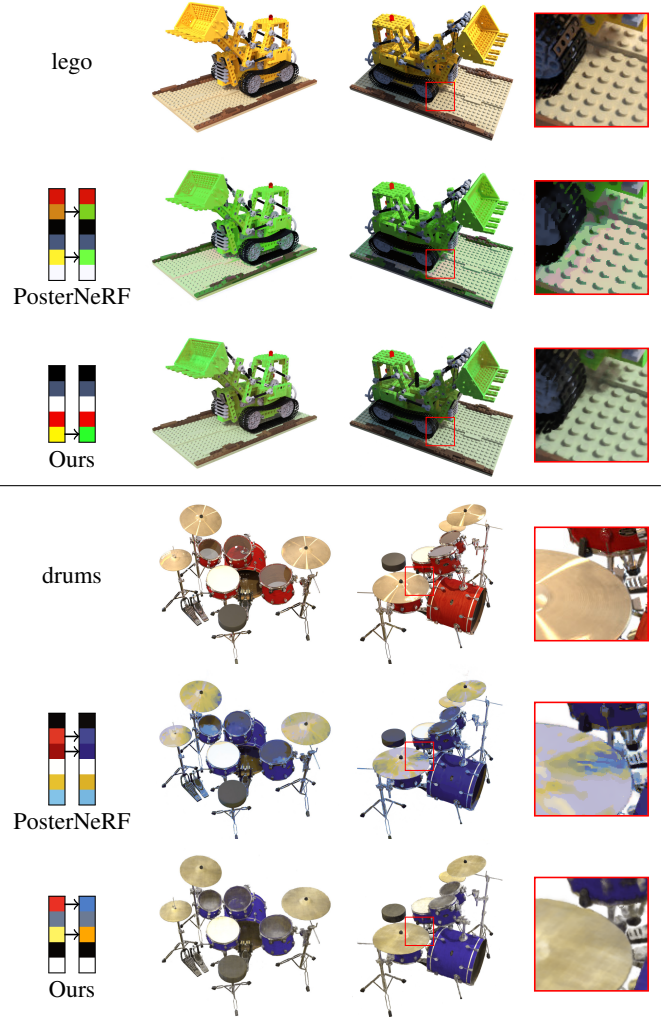


Figure 7: Comparison with PosterNeRF [40]. For each example, from top to bottom, we show the GT original view from dataset, recolored results generated by PosterNeRF and by our PaletteNeRF, respectively. Palettes on the left illustrate the editing operations used by each method.

## 6.2. Comparisons

**Comparison with PosterNeRF.** Figure 7 compares our edited results with a concurrent work – PosterNeRF [40] on two scenes *lego* and *drums*. PosterNeRF also employs a palette-based editing interface. However, it could easily generate visible artifacts, i.e., producing inconsistent recoloring results. Please see the close-up images in Figure 7: PosterNeRF produces unnatural recoloring on the ground near the excavator in scene *lego* and generates unsmooth recoloring on the cymbal in scene *drums*. In contrast, our recoloring results are much better without artifacts.

The reason why PosterNeRF generates artifacts lies in

| method | lego | | | chair | | | flower | | | orchids | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| NeRF | 31.153 | 0.952 | 0.025 | 31.322 | 0.956 | 0.041 | 28.147 | 0.887 | 0.053 | 21.325 | 0.750 | 0.109 |
| PaletteNeRF | 30.135 | 0.941 | 0.030 | 30.723 | 0.950 | 0.045 | 28.172 | 0.881 | 0.059 | 21.515 | 0.759 | 0.108 |

| method | vasedeck | | | valley | | | Cornell box | | | breakfast | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| NeRF | 24.713 | 0.715 | 0.163 | 27.099 | 0.760 | 0.137 | 44.151 | 0.993 | 0.007 | 39.952 | 0.971 | 0.03 |
| PaletteNeRF | 24.657 | 0.721 | 0.152 | 26.638 | 0.745 | 0.159 | 41.339 | 0.985 | 0.013 | 38.875 | 0.967 | 0.035 |

Table 3: Quantitative comparison of PaletteNeRF to the vanilla NeRF on 8 scenes. We report PSNR/SSIM (higher is better) and LPIPS (lower is better). PaletteNeRF is able to achieve similar performance compared to the vanilla NeRF.
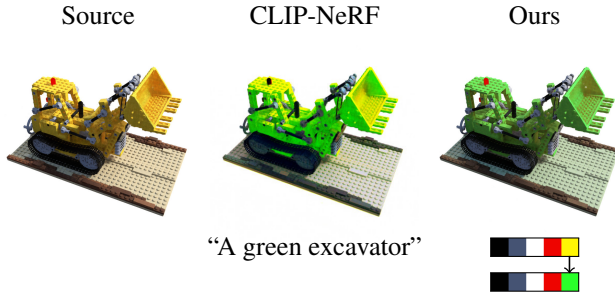


"A green excavator"

Figure 8: Comparison with CLIP-NeRF [42], which takes "A green excavator" as text edit prompt to edit the scene. We change one of the palette colors from yellow to green.
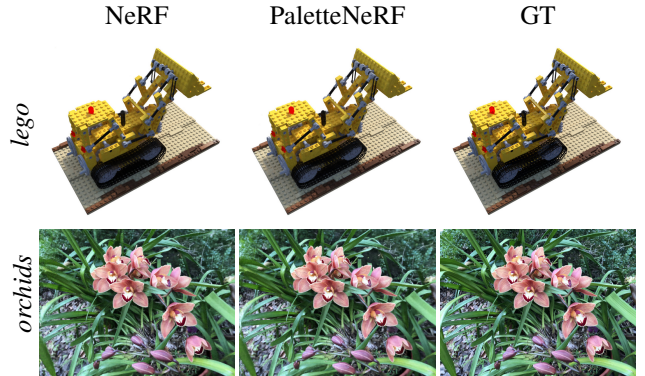


Figure 9: Qualitative comparison of PaletteNeRF to the vanilla NeRF. For each scene, from left to right, we show novel views synthesised by the vanilla NeRF and by our method, respectively, and the ground truth.

(such as the standard ones like *lego* and *drums*).

the way how it uses palettes and computes mixing weights. PosterNeRF will choose at most 2 palette colors, resulting in at most 2 non-zero values of weight $\mathbf{w}^P$ in Eq. 3. Thus some colors can not be reproduced precisely. Furthermore, the weights are also quantized with large step sizes, which leads to color banding artifacts. In contrast, our method approximates each pixel with a linear blending of an arbitrary number of palettes, and the weight of linear blending is accurate, thus producing smooth results.

**Comparison with CLIP-NeRF.** In Figure 8, we further compare our color editing results with CLIP-NeRF [42] on scene *lego*. We can see that CLIP-NeRF tends to blur the detail of the excavator, add noise to the ground, and modify undesired regions such as red lights and grey-blue shafts. In contrast, our method has better image quality and provides better controllability for color editing.

We do not compare with EditingNeRF [20] since it is hard to make a fair comparison. Our method, like vanilla NeRF, could be trained through a single scene and used to recolor the scene. In contrast, EditingNeRF requires a scene dataset with many instances from the same category to enable editing, disabling its ability to edit a single NeRF scene

### 6.3. Results

In Figures 1 and 10, we show the recoloring results on several scenes using our PaletteNeRF. For each example, the recolored results on multiple novel views are provided. PaletteNeRF can produce consistent recoloring results across multiple views. In Table 3, we also provide the quantitative reconstruction scores (including PSNR, SSIM, and LPIPS) of PaletteNeRF for all input scenes. The reconstruction score of the vanilla NeRF is also provided. Figure 9 gives visual comparisons between the reconstructed results of our PaletteNeRF and the vanilla NeRF. It could be found that the results quality of PaletteNeRF is comparable to the vanilla NeRF from both aspects of quantitative measures and visual results. Without reducing visual qualities, our method enables the ability of color editing for NeRFs.
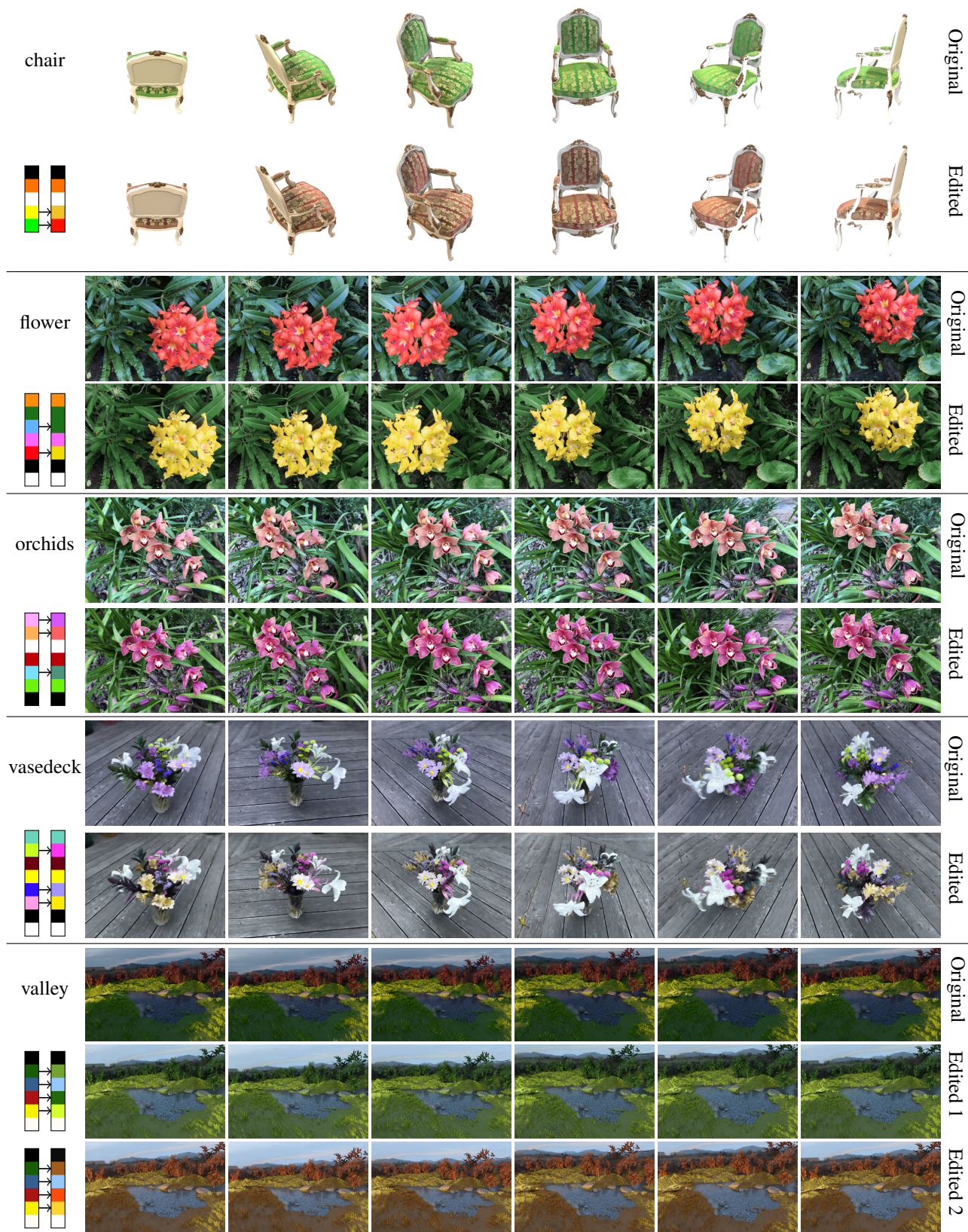
Figure 10: Color editing results of our method. In each 2 rows, the top one is NeRF's output, the bottom one is editing results of PaletteNeRF.

## 7. Conclusion

In this paper, we present PaletteNeRF, which unifies the palette-based image decomposition methods and NeRFs to enable color editing of NeRF-represented scenes. Our method is intuitive, efficient, view-consistent, and artifact-free. The users could recolor the scene by adjusting palette colors, and previewing the recolored results from any novel views. Moreover, the editing process is decoupled from the network structure, which means the backbone of Palette-NeRF and the data preparation step can be replaced with more advanced follow-ups.

Nevertheless, several limitations still exist in our method. For example, we only allow global color editing, hence, if a scene contains two red apples, it is not allowed to only change the color of one apple. Furthermore, geometry editing is not supported. It is worthwhile to investigate ways to edit geometries.

# References

[1] E. Aharoni-Mack, Y. Shambik, and D. Lischinski. Pigment-based recoloring of watercolor paintings. In *Proceedings of the Symposium on Non-Photorealistic Animation and Rendering*, pages 1–11, 2017. 2

[2] Y. Aksoy, T. O. Aydin, A. Smolić, and M. Pollefeys. Unmixing-based soft color segmentation for image manipulation. *ACM Transactions on Graphics (TOG)*, 36(2):1–19, 2017. Publisher: ACM New York, NY, USA. 2

[3] Y. Cao, A. B. Chan, and R. W. Lau. Mining probabilistic color palettes for summarizing color use in artwork collections. In *SIGGRAPH Asia 2017 Symposium on Visualization*, pages 1–8, 2017. 2

[4] H. Chang, O. Fried, Y. Liu, S. DiVerdi, and A. Finkelstein. Palette-based photo recoloring. *ACM Trans. Graph.*, 34(4):139–1, 2015. 2

[5] C.-K. T. Chao, K. Singh, and Y. Gingold. PosterChild: Blend-Aware Artistic Posterization. In *Computer Graphics Forum*, volume 40, pages 87–99. Wiley Online Library, 2021. Issue: 4. 1

[6] F. Dellaert and L. Yen-Chen. Neural Volume Rendering: NeRF And Beyond. *arXiv:2101.05204 [cs]*, Jan. 2021. arXiv: 2101.05204. 2

[7] B. Delos, N. Mellado, D. Vanderhaeghe, and R. Cozot. RGB Point Cloud Manipulation with Triangular Structures for Artistic Image Recoloring. *arXiv preprint arXiv:1912.04583*, 2019. 2

[8] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. CARLA: An open urban driving simulator. In *Conference on robot learning*, pages 1–16. PMLR, 2017. 2

[9] Z.-J. Du, K.-X. Lei, K. Xu, J. Tan, and Y. Gingold. Video recoloring via spatial-temporal geometric palettes. *ACM Transactions on Graphics (TOG)*, 40(4):1–16, 2021. Publisher: ACM New York, NY, USA. 2, 3, 4

[10] Z. Feng, W. Yuan, C. Fu, J. Lei, and M. Song. Finding intrinsic color themes in images with human visual perception. *Neurocomputing*, 273:395–402, 2018. Publisher: Elsevier. 2

[11] M. S. Floater, G. Kós, and M. Reimers. Mean value coordinates in 3D. *Computer Aided Geometric Design*, 22(7):623–631, 2005. Publisher: Elsevier. 3

[12] S. J. Garbin, M. Kowalski, M. Johnson, J. Shotton, and J. Valentin. FastNeRF: High-Fidelity Neural Rendering at 200FPS. pages 14346–14355, 2021. 2

[13] M. Grogan and A. Smolic. Image Decomposition using Geometric Region Colour Unmixing. In *European Conference on Visual Media Production*, pages 1–10, 2020. 2

[14] M. Guo, A. Fathi, J. Wu, and T. Funkhouser. Object-centric neural scene rendering. *arXiv preprint arXiv:2012.08503*, Dec. 2020. 2

[15] T. Ju, S. Schaefer, and J. Warren. Mean value coordinates for closed triangular meshes. In *ACM Siggraph 2005 Papers*, pages 561–566. 2005. 3

[16] S. Kim and S. Choi. Automatic Color Scheme Extraction from Movies. In *Proceedings of the 2020 International Conference on Multimedia Retrieval*, pages 154–163, 2020. 2

[17] S. Lin, M. Fisher, A. Dai, and P. Hanrahan. LayerBuilder: Layer decomposition for interactive image and video color editing. *arXiv preprint arXiv:1701.03754*, 2017. 2

[18] S. Lin and P. Hanrahan. Modeling how people extract color themes from images. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 3101–3110, 2013. 2

[19] L. Liu, J. Gu, K. Z. Lin, T.-S. Chua, and C. Theobalt. Neural Sparse Voxel Fields. In *NeurIPS*, Jan. 2021. arXiv: 2007.11571. 2

[20] S. Liu, X. Zhang, Z. Zhang, R. Zhang, J.-Y. Zhu, and B. Russell. Editing conditional radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5773–5783, 2021. 1, 2, 8

[21] R. Martin-Brualla, N. Radwan, M. S. M. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Jan. 2021. arXiv: 2008.02268. 2

[22] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, Mar. 2020. 1, 2

[23] T. Müller, A. Evans, C. Schied, and A. Keller. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *arXiv:2201.05989 [cs]*, Jan. 2022. arXiv: 2201.05989. 2

[24] R. M. Nguyen, B. Price, S. Cohen, and M. S. Brown. Group-Theme Recoloring for Multi-Image Color Consistency. In *Computer Graphics Forum*, volume 36, pages 83–92. Wiley Online Library, 2017. Issue: 7. 2

[25] M. Niemeyer and A. Geiger. GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. In *CVPR*, page 12, 2021. 2

[26] P. O'Donovan, A. Agarwala, and A. Hertzmann. Color compatibility from large datasets. In *ACM SIGGRAPH 2011 papers*, pages 1–12. 2011. 2

[27] J. Ost, F. Mannan, N. Thuerey, J. Knodt, and F. Heide. Neural scene graphs for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2856–2865, 2021. 2

[28] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R. Martin-Brualla. Nerfies: Deformable Neural Radiance Fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021. 2

[29] K. Park, U. Sinha, P. Hedman, J. T. Barron, S. Bouaziz, D. B. Goldman, R. Martin-Brualla, and S. M. Seitz. HyperNeRF: A Higher-Dimensional Representation for Topologically Varying Neural Radiance Fields. *arXiv:2106.13228 [cs]*, June 2021. arXiv: 2106.13228. 2

[30] M. Pharr, W. Jakob, and G. Humphreys. *Physically based rendering: From theory to implementation*. Morgan Kaufmann, 2016. 6

[31] A. Pumarola, E. Corona, G. Pons-Moll, and F. 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

[32] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, and J. Clark. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 1

[33] C. Reiser, S. Peng, Y. Liao, and A. Geiger. KiloNeRF: Speeding up Neural Radiance Fields with Thousands of Tiny MLPs. In *ICCV 2021*, 2021. 2, 4, 5

[34] J. Tan, S. DiVerdi, J. Lu, and Y. Gingold. Pigmento: Pigment-based image analysis and editing. *Transactions on Visualization and Computer Graphics (TVCG)*, 25(9), 2019. 2

[35] J. Tan, M. Dvorožňák, D. Sykora, and Y. Gingold. Decomposing time-lapse paintings into layers. *ACM Transactions on Graphics (TOG)*, 34(4):1–10, 2015. Publisher: ACM New York, NY, USA. 2

[36] J. Tan, J. Echevarria, and Y. Gingold. Palette-based image decomposition, harmonization, and color transfer. Apr. 2018. 6

[37] J. Tan, J. Echevarria, and Y. Gingold. Efficient palette-based decomposition and recoloring of images via RGBXY-space geometry. *ACM Transactions on Graphics*, 37(6):1–10, Jan. 2019. 2, 3, 4

[38] J. Tan, J.-M. Lien, and Y. Gingold. Decomposing Images into Layers via RGB-Space Geometry. *ACM Transactions on Graphics*, 36(1):1–14, Feb. 2017. 2, 3, 4

[39] A. Tewari, J. Thies, B. Mildenhall, P. Srinivasan, E. Tretschk, Y. Wang, C. Lassner, V. Sitzmann, R. Martin-Brualla, S. Lombardi, T. Simon, C. Theobalt, M. Niessner, J. T. Barron, G. Wetzstein, M. Zollhoefer, and V. Golyanik. Advances in Neural Rendering. *arXiv:2111.05849 [cs]*, Mar. 2022. arXiv: 2111.05849. 2

[40] K. Tojo and N. Umetani. Recolorable Posterization of Volumetric Radiance Fields Using Visibility-Weighted Palette Extraction. In *Computer Graphics Forum*, volume 41, pages 149–160. Wiley Online Library, 2022. Issue: 4. 1, 2, 7

[41] E. Tretschk, A. Tewari, V. Golyanik, M. Zollhofer, C. Lassner, and C. 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

[42] C. Wang, M. Chai, M. He, D. Chen, and J. Liao. CLIP-NeRF: Text-and-Image Driven Manipulation of Neural Radiance Fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. arXiv, Mar. 2022. arXiv:2112.05139 [cs] type: article. 1, 2, 8

[43] Y. Wang, Y. Liu, and K. Xu. An Improved Geometric Approach for Palette-based Image Decomposition and Recoloring. In *Computer Graphics Forum*, volume 38, pages 11–22. Wiley Online Library, 2019. Issue: 7. 2, 3, 4

[44] B. Yang, Y. Zhang, Y. Xu, Y. Li, H. Zhou, H. Bao, G. Zhang, and Z. Cui. Learning object-compositional neural radiance field for editable scene rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13779–13788, 2021. 2

[45] A. Yu, R. Li, M. Tancik, H. Li, R. Ng, and A. Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. 2

[46] H.-X. Yu, L. J. Guibas, and J. Wu. Unsupervised discovery of object radiance fields. In *International Conference on Learning Representations*, 2022. 2

[47] K. Zhang, G. Riegler, N. Snavely, and V. Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2

[48] Q. Zhang, C. Xiao, H. Sun, and F. Tang. Palette-based image recoloring using color decomposition optimization. *IEEE Transactions on Image Processing*, 26(4):1952–1964, 2017. Publisher: IEEE. 2