

# Transforming Radiance Field with Lipschitz Network for Photorealistic 3D Scene Stylization

Zicheng Zhang<sup>1</sup> Yinglu Liu<sup>2</sup> Congying Han<sup>1\*</sup> Yingwei Pan<sup>2</sup> Tiande Guo<sup>1</sup> Ting Yao<sup>2</sup>  
<sup>1</sup>University of Chinese Academy of Sciences <sup>2</sup>JD AI Research

zhangzicheng19@mails.ucas.ac.cn liuyinglu1@jd.com hancy@ucas.ac.cn  
panyw.ustc@gmail.com tdguo@ucas.ac.cn tingyao.ustc@gmail.com

## Abstract

Recent advances in 3D scene representation and novel view synthesis have witnessed the rise of Neural Radiance Fields (NeRFs). Nevertheless, it is not trivial to exploit NeRF for the photorealistic 3D scene stylization task, which aims to generate visually consistent and photorealistic stylized scenes from novel views. Simply coupling NeRF with photorealistic style transfer (PST) will result in cross-view inconsistency and degradation of stylized view syntheses. Through a thorough analysis, we demonstrate that this non-trivial task can be simplified in a new light: **When transforming the appearance representation of a pre-trained NeRF with Lipschitz mapping, the consistency and photorealism across source views will be seamlessly encoded into the syntheses.** That motivates us to build a concise and flexible learning framework namely LipRF, which upgrades arbitrary 2D PST methods with Lipschitz mapping tailored for the 3D scene. Technically, LipRF first pre-trains a radiance field to reconstruct the 3D scene, and then emulates the style on each view by 2D PST as the prior to learn a Lipschitz network to stylize the pre-trained appearance. In view of that Lipschitz condition highly impacts the expressivity of the neural network, we devise an adaptive regularization to balance the reconstruction and stylization. A gradual gradient aggregation strategy is further introduced to optimize LipRF in a cost-efficient manner. We conduct extensive experiments to show the high quality and robust performance of LipRF on both photorealistic 3D stylization and object appearance editing.

## 1. Introduction

Photorealistic style transfer (PST) [52] is one of the important tasks for visual content creation, which aims to automatically apply the color style of a reference image to an

other input (e.g., image [38] or video [65]). In this task, the stylized result is required to look like a camera shot and preserve the input structure (e.g., edges and regions). Benefiting from the launch of deep learning, a series of sophisticated deep PST methods [19, 30, 38, 59, 65] have been developed for practical usage. Recent progress in 3D scene representation has featured Neural Radiance Field [43, 67] (NeRF) with efficient training and high-quality view synthesis. This inspires us to go one step further to explore a more challenging task of photorealistic 3D scene stylization, which is to generate visually consistent and photorealistic stylized syntheses from arbitrary views. Such a task enables an automatic modification of 3D scene appearance with different lighting, time of day, weather, or other effects, thereby enhancing user experience and stimulating emotions for virtual reality [57].

Nevertheless, it is not trivial to build an effective framework for photorealistic 3D scene stylization. The difficulty mainly originates from the fact that there is no valid photorealistic style loss tailored for training NeRF. In general, the image-based PST is commonly tackled via either the neural style transfer [11] combined with complicated post-processing [30, 38, 41], or particular network structures [29, 61, 65]. However, none of them can be directly applied to the learning of NeRF. As shown in Figure 1, simply employing the state-of-the-art 2D PST on each view might result in noise, disharmony and even inconsistency across views, since the PST methods rely on the size or object masks of the inputs. Such downsides will be further amplified after reconstructing the 3D scene with NeRF.

To alleviate these limitations, we start with a basic understanding of this task: *Though preserving the photorealism and consistency seems to be different in the context of 2D images, they do have the same essence when moving to 3D volume rendering [40].* From this standpoint, the task is simplified as a problem to regulate the volume rendering variance of the radiance field before and after stylization. According to the studies of color mapping [48, 51, 52], some

\*Corresponding author

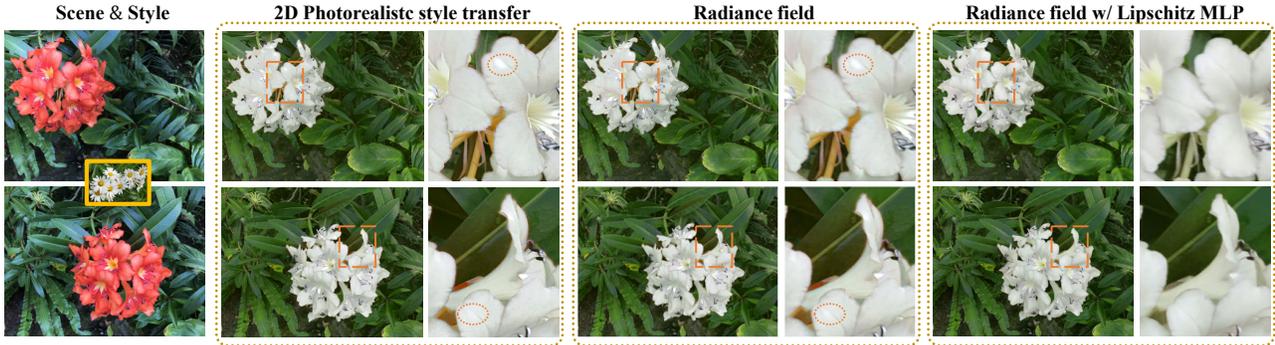


Figure 1. Illustrations of different strategies for photorealistic 3D scene stylization. The 2D PST method [9] generates disharmonious orange color on the stem and petaline edge. The regions bounded by the dotted ellipses are inconsistent due to the white spot in the first view. When employing a radiance field to reconstruct the results of 2D PST method, the disharmony and inconsistency are still retained. Our LipRF successfully eliminates these downsides, and renders high-quality stylized view syntheses.

specific linear mappings of image pixels can nicely preserve the image structures with photorealistic effect. Motivated by this, we theoretically demonstrate that a simple yet effective design of Lipschitz-constrained linear mapping over appearance representation can elegantly control the volume rendering variance. Furthermore, we prove that replacing linear mapping with a Lipschitz multilayer perceptron (MLP) also holds these properties under extra assumptions which can be relaxed in practice. Such a way completely eliminates the drawbacks of 2D PST when transforming the radiance field with a Lipschitz MLP (see Figure 1). In a nutshell, our analysis verifies that Lipschitz MLP can be interpreted as an implicit regularization to safeguard the 3D photography of stylized scenes.

By consolidating the idea of transforming the radiance field with the Lipschitz network, we propose a novel NeRF-based architecture (namely LipRF) for photorealistic 3D scene stylization. Technically, LipRF contains two stages: 1) training a radiance field to reconstruct the source 3D scene; 2) learning a Lipschitz network to transform the pre-trained appearance representation to the stylized 3D scene with the guidance of style emulation on each view by arbitrary 2D PST. We adopt the Plenoxels [67] as the base radiance field due to its advanced reconstruction quality and compressed appearance representation by spherical harmonics. Considering that the Lipschitz condition greatly impacts the expressivity of neural networks, we design an adaptive regularization based on spectral normalization [44] to allow a mild relaxation of the Lipschitz constant in each linear layer. Finally, we capitalize on gradual gradient aggregation to optimize LipRF in a cost-efficient fashion.

In summary, we have made the following contributions: **(I)** We present a thorough and insightful analysis of photorealistic 3D scene stylization on the basis of volume rendering and Lipschitz transformation. **(II)** We build a concise and flexible framework (LipRF) for photorealistic 3D scene stylization by novelly transforming the appearance representation of a pre-trained NeRF with the Lipschitz Network. **(III)** Under the Lipschitz condition, we design adaptive regularization and gradual gradient aggregation to seek a better

trade-off among the reconstruction, stylization quality, and computational cost. We evaluate LipRF on both photorealistic 3D stylization and object appearance editing tasks to validate the effectiveness of our proposal.

## 2. Related works

### 2.1. Novel view synthesis

Novel view synthesis aims to synthesize view images with arbitrary camera poses from given source images. Many works have been proposed on various discrete representations, *e.g.*, multi-plane image [42, 72], point clouds [13, 32], meshes [10, 58], and voxels [35]. Recently, neural radiance field approaches [4, 39, 43, 70] encode the scene into a continuous implicit volumetric representation via multilayer perceptron, and render the novel view by volume rendering integral [40]. Later, Plenoxels [67] reveals that the key element of NeRF is the differentiable volume renderer. By simplifying NeRF into a sparse voxel grid with spherical harmonics, Plenoxels achieves substantial speedups with comparable rendering quality. Our work capitalizes on the benefits of Plenoxels, and enables photorealistic 3D scene stylization to achieve in a few minutes.

### 2.2. Style transfer methods

**Image stylization.** Photorealistic style transfer is a long-standing topic [2, 5, 48, 50, 52] that focuses on color distribution transfer while maintaining the photorealism of the image. In the regime of deep learning, since Gatys *et al.* [11] present the neural style transfer that matches the feature distribution [6, 21, 24, 31, 49, 71] to transfer artistic texture, many works incorporate neural networks into PST for more complicated visual effects. One of the early works, Luan *et al.* [38] introduces the semantic mask and photorealistic regularization [26] to yield impressive results. PhotoWCT [30] further improves the regularization [26] into a closed form. The works [9, 27, 41, 65] propose to preserve the high frequency to ensure photorealism. For example, WCT<sup>2</sup> [65] embeds the wavelet transform into neural networks to fully retain image structure. Other works adopt

particular style transfer operators [17,19,28,29,37,59] to the Encoder-Decoder pre-trained on natural image dataset, *e.g.*, COCO [33]. The advances [60,61] learn the edge-preserved local affine grid [12] in bilateral space [5] to transfer the color locally. Because the input size and semantic mask heavily influence these methods, directly applying to 3D scene will produce distortion and inconsistency.

**3D scene stylization.** Recently, several works [8,18,47,69] couple NeRF with neural style transfer [11] to tackle artistic 3D scene stylization. For example, [8,47,69] finetune pre-trained NeRFs with the image-based style losses. [18] adapts the pre-trained NeRF by means of a 2D style transfer model. These methods enable better stylization in a variety of scenes, and perform faster training and inference than the approaches [16,46,64]. However, strict consistency and photorealism are hardly achieved by these methods, making them inapplicable to photorealistic stylization.

### 2.3. Lipschitz network

The Lipschitz network [54], *i.e.*, neural network with a limited Lipschitz constant, has advantages in robustness [63], generalization [66], and stability of training [44]. Since computing the exact Lipschitz constant of neural networks is NP-hard [54], previous methods usually minimize [14,34] or specify [15,44] an upper bound of Lipschitz constant as a small value. Due to the small Lipschitz constant limits the visual effect in synthesis task, these methods cannot provide a proper constraint for the renderer/generator. To solve this, we propose to adaptively constrain Lipschitz property according to the given scene and reference image.

## 3. Preliminaries

**Radiance field.** Generally, the radiance field [43] is a 5D function  $\mathbf{F}$  that maps any 3D location  $\mathbf{x}$  and viewing direction  $\mathbf{d}$  to volume density  $\sigma$  and color  $\mathbf{c} = (r, g, b)$ . Specifically,  $\mathbf{F}$  can be further divided into the geometry part  $\sigma = \mathbf{F}_{geo}(\mathbf{x})$  and the appearance part  $\mathbf{c} = \mathbf{F}_{app}(\mathbf{x}, \mathbf{d})$ , *i.e.*,  $\mathbf{F} = (\mathbf{F}_{geo}, \mathbf{F}_{app})$ . Recently, [67,68] factorize the appearance via spherical harmonic representation:

$$\mathbf{c} = \mathbf{F}_{sh}(\mathbf{x})\Gamma(\mathbf{d}) + \mathbf{v}, \quad (1)$$

where  $\Gamma$  is the pre-defined basis function to produce the  $\ell$ -dimensional basis,  $\mathbf{F}_{sh}$  computes the corresponding  $3 \times \ell$  coefficient matrix, and  $\mathbf{v}$  is a fixed vector to normalize the colors. This form can greatly reduce the redundancy of the representation, and speed up the training and inference.

**Volume rendering.** A ray cast into the scene can be formulated as  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ , where  $\mathbf{o}$  and  $\mathbf{d}$  are the origin and normalized direction of the ray, and  $t$  denotes the distance along the ray. The color of the ray with near and far bounds  $t_1$  to  $t_{T+1}$  is estimated by the volume rendering

$$C(\mathbf{r}; \mathbf{F}) = \sum_{i=1}^T w_i \mathbf{c}_i, \quad (2)$$

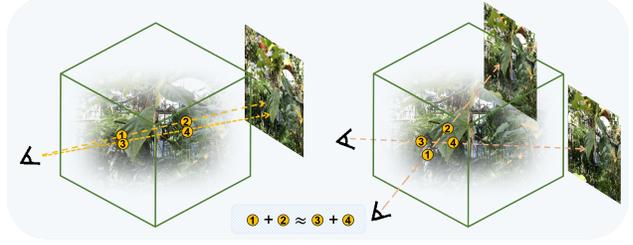


Figure 2. Illustrations of image structure (left) and cross-view consistency (right) from the perspective of volume rendering. For easy understanding, we placed images behind radiance fields (the cubes), and the two yellow circles on each ray denote the weighted colors sampled for volume rendering as Eq. (2). For each case, the volume rendering variance as Eq. (5) should be a small value.

$$w_i = (1 - e^{-\sigma_i(t_{i+1}-t_i)})e^{-\sum_{i' < i} \sigma_{i'}(t_{i'+1}-t_{i'})}, \quad (3)$$

where  $\sigma$  and  $\mathbf{c}$  are predicted by the radiance field  $\mathbf{F}$ .

**Lipschitz functions.** Given two metric spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is Lipschitz, or  $K$ -Lipschitz continuous if there exists  $K \in \mathbb{R}^+$  satisfied that:

$$\forall x_1, x_2 \in \mathcal{X}, d_{\mathcal{Y}}(f(x_1), f(x_2)) \leq K d_{\mathcal{X}}(x_1, x_2), \quad (4)$$

where  $d_{\mathcal{X}}$  and  $d_{\mathcal{Y}}$  are metrics in  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively.  $K$  is called the Lipschitz constant of  $f$ . In this paper, we define  $d_{\mathcal{X}}$  and  $d_{\mathcal{Y}}$  as (2-norm) Euclidean distances of vectors.

## 4. Problem statement

**Task and challenges.** Given a set of images  $\{\mathbf{I}_i\}_{i=1}^N$  taken in a 3D scene  $\mathcal{I}$  with known camera parameters, our goal is to synthesize novel photorealistic views with the similar color style to the reference image. As discussed in Section 1, the main challenge is to protect *photorealism* and *consistency* while transferring the color style. More specifically, 1) For photorealism, the image structure (*e.g.*, edges and regions) needs to be well preserved after stylization. Without loss of generality, considering two nearby pixels at  $\mathbf{p}_1$  and  $\mathbf{p}_2$  in an image  $\mathbf{I}$ , we succinctly state that if  $\|\mathbf{I}(\mathbf{p}_1) - \mathbf{I}(\mathbf{p}_2)\| < \epsilon$  ( $\epsilon$  is a small value), the two belong to the same region; otherwise not. 2) For consistency, we suppose that one point in the scene can be observed by two adjacent views  $\mathbf{I}_1$  and  $\mathbf{I}_2$  at  $\mathbf{p}_1$  and  $\mathbf{p}_2$  (with a slight abuse of notation), respectively. Thus, the cross-view consistency should meet  $\|\mathbf{I}_1(\mathbf{p}_1) - \mathbf{I}_2(\mathbf{p}_2)\| < \epsilon$ . For this task, the above two relations between pixels in the stylized images should be consistent with that of the source images.

**A novel perspective from volume rendering.** Due to the fact that the intensity of pixels is physically derived from the volume rendering of corresponding ray castings, we reinterpret them from the perspective of rays. Suppose that  $\mathbf{r}_1$  and  $\mathbf{r}_2$  intersect with the single/couple views at  $\mathbf{p}_1$  and  $\mathbf{p}_2$ , referring to Eq. (2), we define the volume rendering variance

( $vr_r$ ) of rays  $\mathbf{r}_1$  and  $\mathbf{r}_2$  in the radiance field  $\mathbf{F}$  by

$$\begin{aligned} vr_r(\mathbf{r}_1, \mathbf{r}_2; \mathbf{F}) &= \|C(\mathbf{r}_1; \mathbf{F}) - C(\mathbf{r}_2; \mathbf{F})\| \\ &= \left\| \sum_{i=1}^T w_i^{r_1} \mathbf{c}_i^{r_1} - \sum_{i=1}^T w_i^{r_2} \mathbf{c}_i^{r_2} \right\|, \end{aligned} \quad (5)$$

where the superscripts denote the ray index. Importantly, the arrival of  $vr_r$  integrates the intricate relationships within (structure) and between (consistency) images into the variance over rays. Figure 2 depicts a concise example. In this way, if a radiance field  $\mathbf{F}$  could represent the scene  $\mathcal{I}$ , the core challenge of this task can be streamlined to learn a stylized radiance field  $\mathbf{F}'$  satisfying

$$vr_r(\mathbf{r}_1, \mathbf{r}_2; \mathbf{F}) < \epsilon \Rightarrow vr_r(\mathbf{r}_1, \mathbf{r}_2; \mathbf{F}') < \epsilon', \quad (6)$$

where  $\epsilon'$  is a small value. In the following, we prove that the above demand can be fulfilled elegantly recurring to the Lipschitz mapping, thereby leading to decent stylization.

## 5. Methodology

In this section, we introduce a concise and flexible framework called LipRF to tackle photorealistic 3D scene stylization. LipRF first obtains the radiance field  $\mathbf{F}$  of the source scene (Sec. 5.1). Based on the theoretical analysis (Sec. 5.2) of controlling  $vr_r$  with Lipschitz mappings, LipRF transforms the pre-trained radiance field with the Lipschitz MLP (Sec. 5.3) to reconstruct the views stylized by 2D PST. The gradual gradient aggregation (Sec. 5.4) is elaborated to optimize LipRF in a cost-efficient way.

### 5.1. Scene representation via radiance field

In the first stage, a radiance field  $\mathbf{F} = (\mathbf{F}_{geo}, \mathbf{F}_{app})$  with faithful geometry and appearance representation is trained to reconstruct the real scene  $\mathcal{I}$ . Following [67], we adopt the reconstruction loss for training:

$$\mathcal{L}_{rec}(\mathbf{F}, \mathcal{I}) = \sum_{i=1}^m \|C(\mathbf{r}_i; \mathbf{F}) - C(\mathbf{r}_i)\|^2, \quad (7)$$

where  $\{\mathbf{r}_i\}_{i=1}^m$  denotes the set of rays generated under the given camera parameters,  $C(\mathbf{r}_i; \mathbf{F})$  is the color estimated by volume rendering as in Eq. (2), and  $C(\mathbf{r}_i)$  is the groundtruth color of the corresponding image pixel in  $\{\mathbf{I}_i\}_{i=1}^N$ .

We assume  $\mathbf{F}_{geo}$  enables fully encoding the geometry of  $\mathcal{I}$  after training. Since PST should not change the geometry of the source scene, we directly set  $\mathbf{F}'_{geo} = \mathbf{F}_{geo}$ . In this way,  $C(\mathbf{r}; \mathbf{F})$  and  $C(\mathbf{r}; \mathbf{F}')$  have the same rendering weights as in Eq. (3) on the ray path.

### 5.2. Theoretic form of stylized radiance field

The classic 2D PST methods [51, 52, 62], which simply transfer the color style by linear mappings of pixels, can well preserve the image structure. In addition to linearity, we find another commonality of these methods that the

corresponding Lipschitz constants are all of small values, e.g., usually less than 5 on the PST dataset [38]. The fact indicates that Lipschitz property may also play an important role in maintaining structure of images. Prompted by this underlying relationship, we prove that the Lipschitz-constrained linear mapping of  $\mathbf{F}_{app}$  is indeed an optimal form for holding Cond. (6):

**Proposition 1.** *Considering  $f(\mathbf{c}) = \mathbf{A}\mathbf{c} + \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{3 \times 3}$ ,  $\mathbf{b} \in \mathbb{R}^{3 \times 1}$ , if  $\mathbf{F}'_{app} = f \circ \mathbf{F}_{app}$ ,  $\sum_{i=1}^T w_i = 1$  and  $vr_r(\mathbf{r}_1, \mathbf{r}_2; \mathbf{F}) < \epsilon$ , we have  $vr_r(\mathbf{r}_1, \mathbf{r}_2; \mathbf{F}') < K\epsilon$ , where  $K = \|\mathbf{A}\|_2$  is the Lipschitz constant of  $f^1$ .*

Since the unique assumption of  $\sum_{i=1}^T w_i = 1$  is in accordance with many radiance field models [43, 67, 68] in practice, Prop. 1 provides an ideal way to establish  $\mathbf{F}'$ , where the variational bound  $\epsilon'$  in Cond. (6) is influenced by the Lipschitz constant of transform. Nonetheless, the limited expressivity of linear mapping definitely affects dramatic style effects, e.g., the color variation in Figure 1 by the deep PST method. We propose to mitigate the deficiency while maintaining Cond. (6) by means of the Lipschitz MLP:

**Proposition 2.** *Considering  $f = f_l \circ \dots \circ f_1$ ,  $f_j(x) = \mathbf{A}_j x + \mathbf{b}$  if  $j = l$  and  $\sigma(\mathbf{A}_j x)$  otherwise, where  $\sigma = \max(0, x)$ . If  $\mathbf{F}'_{app} = f \circ \mathbf{F}_{app}$ ,  $\sum_{i=1}^T w_i = 1$  and  $\max_{i=1, \dots, T} \|w_i^{r_1} \mathbf{c}_i^{r_1} - w_i^{r_2} \mathbf{c}_i^{r_2}\| < \epsilon/T$ , we have  $vr_r(\mathbf{r}_1, \mathbf{r}_2; \mathbf{F}') < K\epsilon$ , where  $K = \prod_{i=1}^l \|\mathbf{A}_i\|_2$  is the Lipschitz constant of  $f^1$ .*

This proposition further necessitates the vanishing of  $\|w_i^{r_1} \mathbf{c}_i^{r_1} - w_i^{r_2} \mathbf{c}_i^{r_2}\|$ , which is valid when the rays are quite adjacent. Despite the premise seems to hamper the utilization of Lipschitz MLP, it can be greatly loosened for the close relation between Lipschitz MLP and linear mapping in Prop. 1: 1) The above Lipschitz MLP as a piece-wise linear function behaves the same as linear mapping in a local space [45]. 2) The Lipschitz MLP with strict constraint of Lipschitz condition and gradient norm approximates to the linear mapping [1]. These properties encourage to form the stylized radiance field as:

$$\mathbf{F}' = (\mathbf{F}_{geo}, f \circ \mathbf{F}_{app}), \quad f \text{ is } K\text{-Lipschitz MLP}. \quad (8)$$

Regarding the training complexity of  $f$ , generally if there are  $n$  values sampled for each variable of the position and direction,  $f$  needs to take a large amount of parameters with high computational costs for predicting  $O(n^5)$  colors correctly. Fortunately, we find a nice property<sup>1</sup> in Eq. (1):

$$\mathbf{A}\mathbf{F}_{app}(\mathbf{x}, \mathbf{d}) + \mathbf{b} \Leftrightarrow \mathbf{A}\mathbf{F}_{sh}(\mathbf{x}) + 2\sqrt{\pi}[\mathbf{A}\mathbf{v} + \mathbf{b} - \mathbf{v}, \mathbf{0}], \quad (9)$$

which allows to exchange the linear mappings of appearance representation and spherical harmonic representation.

<sup>1</sup> Proof is provided in the supplementary materials.

Therefore,  $\mathbf{F}'_{app} = f \circ \mathbf{F}_{app}$  can be pared down to  $\mathbf{F}'_{sh} = f \circ \mathbf{F}_{sh}$ , while not violating the above propositions. By doing so, it is feasible to design  $f$  as a lightweight model to handle the  $O(n^3)$  spherical harmonic coefficients.

### 5.3. Lipschitz transformation of radiance field

Based on the above analysis, the second step of LipRF is to transform the radiance field  $\mathbf{F}$  with Lipschitz MLP  $f$  as stylized radiance field  $\mathbf{F}'$ . Here  $f$  is composed of linear and activation layers. It receives and updates the flattened spheric harmonic coefficients. We also input the 3D position for spatial inductive bias, namely  $\mathbf{F}'_{sh}(\mathbf{x}) = f(\mathbf{F}_{sh}(\mathbf{x}), \mathbf{x})$ . The training objective is

$$\min_f \mathcal{L}_{rec}(\mathbf{F}', \mathcal{S}) + \lambda \mathcal{L}_{Lip}(f). \quad (10)$$

$\mathcal{S}$  denotes the scene consisting of  $\{pst(\mathbf{I}_i)\}_{i=1}^N$ , where  $pst$  is an arbitrary 2D PST method.  $\mathcal{L}_{rec}(\mathbf{F}', \mathcal{S})$  is the reconstruction loss of the same form as Eq. (7), and  $\mathcal{L}_{Lip}$  is the proposed adaptive Lipschitz regularization to adjust the Lipschitz constant of  $f$ , and  $\lambda$  is the balance weight.

**Lower bound of Lipschitz constant.** In practice, it is difficult to determine the optimal Lipschitz constant  $K$  of the network for different scenes and reference images. Large values will invalidate the constraint of  $vr_r$ , while small values may result in over-constrained conditions for the color style transfer. To address this problem, we first estimate a value  $K_{est}$  as the lower bound of  $K$ . Here we adopt a widely used color transfer method, *i.e.*, Monge-Kantorovitch linear [51] (MKL) mapping, to compute the transfer matrix between  $\mathcal{I}$  and  $\mathcal{S}$ :

$$\mathbf{M} = \Sigma_{\mathcal{I}}^{-1/2} \left( \Sigma_{\mathcal{I}}^{1/2} \Sigma_{\mathcal{S}} \Sigma_{\mathcal{I}}^{1/2} \right)^{1/2} \Sigma_{\mathcal{I}}^{-1/2}, \quad (11)$$

where  $\Sigma_{\mathcal{I}}$  and  $\Sigma_{\mathcal{S}}$  are the covariance matrices of pixel colors in  $\{\mathbf{I}_i\}_{i=1}^N$  and  $\{pst(\mathbf{I}_i)\}_{i=1}^N$ , respectively. We specify  $K_{est} = \|\mathbf{M}\|_2$  so that  $f$  outperforms the linear mapping.

**Adaptive Lipschitz regularization.** Since the Lipschitz constant of network is affected by that of each linear layer (see Prop. 2), it is hard to optimize  $K$  directly. We reform

$$\mathbf{A}_i = \text{squareplus}(K_i, b) \mathbf{W}_i / \|\mathbf{W}_i\|_2, \quad (12)$$

where  $\mathbf{W}_i$  and  $K_i$  are the parameters to optimize for the  $i$ -th linear layer. Here,  $\text{squareplus}(x, b) = \frac{1}{2} (x + \sqrt{x^2 + b})$  [3] is similar to ReLU, but always produces positive values to prevent the norm of  $\mathbf{A}_i$  from vanishing. We follow [44] to fast approximate  $\|\mathbf{W}_i\|_2$  by one-step power iteration. On this basis, the regularization is defined as

$$\mathcal{L}_{Lip} = \sum_{i=1}^l \text{squareplus}(K_i - \sqrt{K_{est}}, b). \quad (13)$$

The Lipschitz constant of network is optimized by constraining the norm of each linear layer to the geometric mean value of  $K_{est}$ , and thus  $\mathcal{L}_{Lip}$  can softly control the gap between  $K$  and  $K_{est}$  during training.

---

### Algorithm 1: PyTorch-style pseudocode for GGA

---

```
# f - Lipschitz MLP; Opt - optimizer of f; sh -
[B,ℓ] spheric harmonic coefficients of F; sigma
- [B,1] density of F; rs - rays of a view; C -
groundtruth; idx - indexes for splitting batches
Opt.zero_grad()
# 0. feed forward
With torch.no_grad():
    sh_t = cat([f(sh[i:j]) for i,j in idx])
    C_hat = volume_render(rs, sigma, sh_t)
# 1. backward from rec loss to image
rec_loss(C_hat, C).backward()
# 2. backward from image to sh_t
p = volume_render(rs, sigma, sh_t)
p.backward(grad = C_hat.grad)
# 3. backward from sh_t and Lip loss to f
for i, j in idx:
    p = f(sh[i:j])
    p.backward(grad = sh_t[i:j].grad/B)
(lambda * Lip_loss(f)).backward()
Opt.step()
```

---

### 5.4. Optimization by gradual gradient aggregation

For Plenoxels [67], Lipschitz MLP needs to transform the spherical harmonic coefficients on all vertices at each iteration due to the structure of voxel grid. Therefore, training LipRF will cost a massive GPU memory in both feed forward and back propagation. The intuitive strategy that optimizes a sparse set of rays [8,47] at once will cause considerable redundancy, since the majority of vertices are not selected. [69] proposes to defer the back propagation for removing useless gradient cache, but it still fails on LipRF due to the huge amount of inputs for Lipschitz MLP.

To increase training efficiency, we propose the gradual gradient aggregation (GGA) detailed in Algorithm 1. First, GGA does not construct the computation graphs during the forward process to reduce memory footprint. Then, the GGA gradually propagates the gradient after redoing each forward step. Finally, the gradient of Lipschitz MLP is aggregated in a batch-wise way. Since the forward process costs much less time than the backward propagation, GGA enables a fast training speed, and tractable memory footprint by adjusting the number of batches.

## 6. Experiments

**Settings.** The overall architecture of LipRF is implemented based on the official code of Plenoxels<sup>2</sup> [67]. The first stage of training radiance field is the same as in Plenoxels. In the second stage of learning Lipschitz network, the MLP has 5 linear-activation layers and 1 linear output layer. The middle layer has 64 neural units. Note that taking the (1-Lipschitz) sinusoidal function [55] as activation can lead to better results than ReLU or LeakyReLU sometimes. We use the Adam optimizer [22], where  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,

<sup>2</sup><https://github.com/sxyu/svox2>

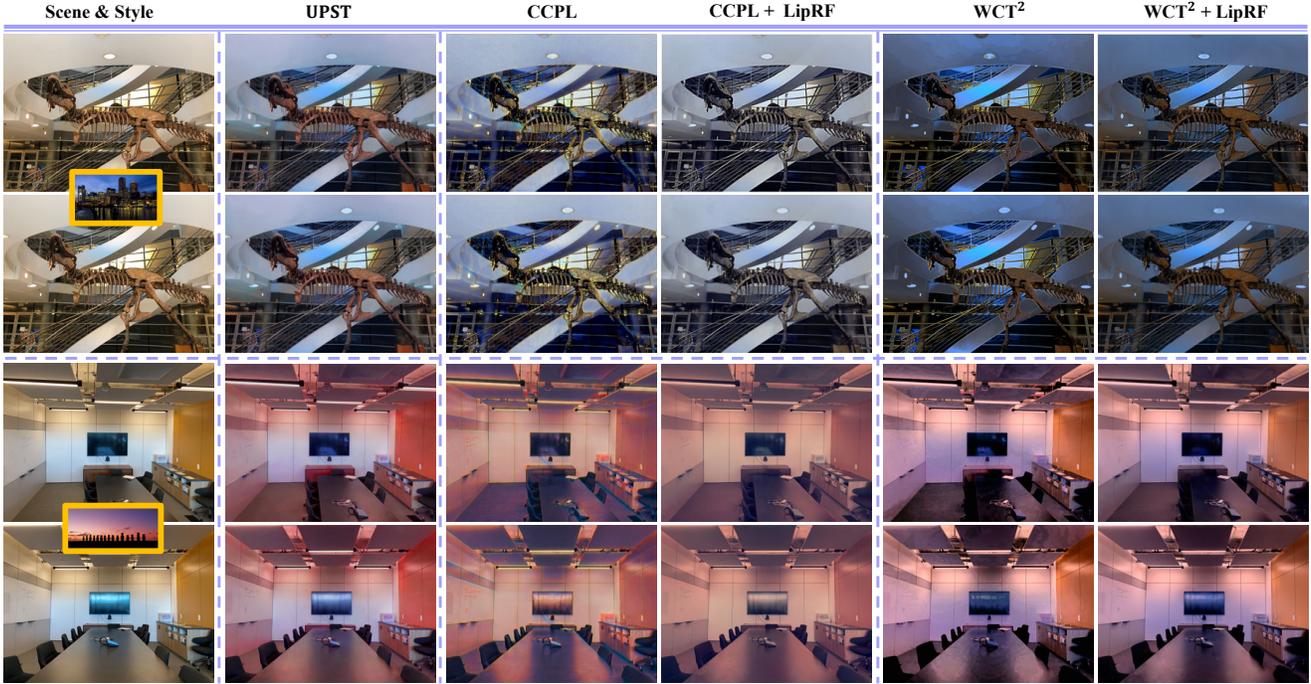


Figure 3. Comparison on “trex” and “room” scenes from LLFF dataset [42], where the image resolution is  $1008 \times 756$ . For the two scenes, LipRF is learned with the prior knowledge derived from the stylized results of  $WCT^2$  [51] and CCPL [51].

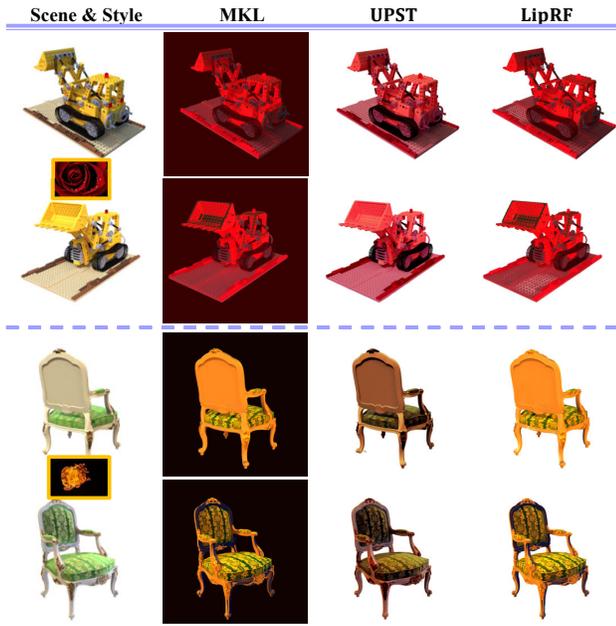


Figure 4. Comparison on NeRF-synthetic dataset [43], where the image resolution is  $800 \times 800$ . LipRF is learned with the prior knowledge derived from the stylized scene of MKL [51].

and the learning rate is reduced from  $10^{-2}$  to  $10^{-4}$  by cosine annealing [36]. We set  $\lambda = 2 \times 10^{-4}$  in Eq. (10) and  $b = 10^{-12}$  in squareplus. Unless otherwise stated, the opti-

mization process runs 300 epochs in total and takes no more than 7 minutes on a single NVIDIA RTX 3090.

**Datasets.** We conduct qualitative and quantitative evaluations to verify LipRF on various scenes in multiple datasets, including *NeRF-synthetic* dataset [43] of synthetic scenes, *LLFF* [42] dataset of real forward-facing scenes, and some scenes from *Tanks and Temples* dataset [23] of real  $360^\circ$  scenes and the multi-view stereo *DTU* dataset [20]. The style images are derived from the PST dataset [38]. Due to space limitations, we present part of them in the paper. Please refer to the supplementary materials for more results.

**Baselines.** We leverage three existing 2D PST approaches as baselines. In particular, MKL [51] is a classic and widely used color transfer method.  $WCT^2$  [65] is recognized as a typical baseline of 2D PST methods, which is stable for videos or high-resolution images. CCPL [59] is a recent state-of-the-art 2D PST method. Besides, we include a concurrent work UPST [7] for comparison, which promotes 2D PST to tackle the same task of photorealistic 3D stylization.

## 6.1. Qualitative results

**NeRF-synthetic dataset.** Figure 4 summarizes the comparison between our LipRF and existing 2D/3D PST methods. Considering that the scenes in synthetic dataset [43] are constructed with simple texture, we choose the classic MKL [51] as the only 2D PST baseline. Although MKL can transfer the colors faithfully, it fails to distinguish between the foreground and background, resulting in a drastic

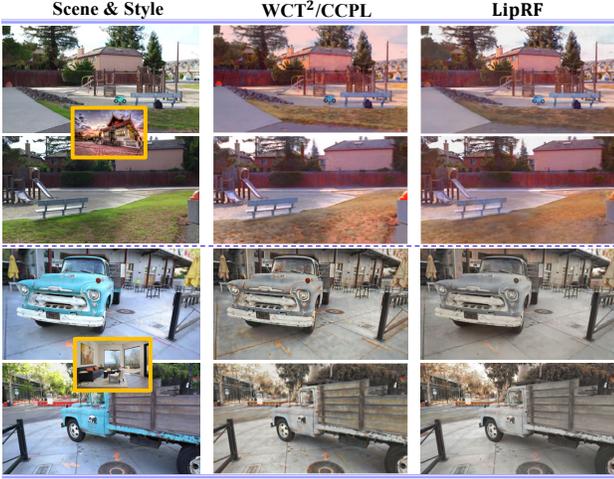


Figure 5. Comparison on Tanks and Temples dataset [23]. For two scenes, LipRF is learned with the prior knowledge derived from the stylized results of WCT<sup>2</sup> (first) [51] and CCPL [51] (second).

change in the background pixel colors. It is worth noting that this is indeed a common problem existing in all 2D PST methods when being simply applied to 3D scene. Moreover, the stylized results of UPST [7] show some disparity with the reference color style in vision. In contrast, our approach manages to integrate the advantages of radiance field with MKL. Since the densities of background regions are all 0 in the radiance field, they do not contribute to the volume rendering, encouraging the rendered colors to be consistent with the source scene background.

**LLFF dataset.** We further illustrate the comparisons on LLFF dataset in Figure 3. As shown in the results, the state-of-the-art 2D PST approach (CCPL [59]) has two main downsides: First, there are a lot of noises in the stylized scene, like contrasting spots and variegation. Second, the stylized results are obviously blurred compared to the source images. This is mainly because it is non-trivial to apply the Encoder-Decoder architecture in CCPL trained on COCO dataset [33] for the high-resolution inputs of scene images. WCT<sup>2</sup> [59] performs better than CCPL in structure preservation due to its wavelet module, but the inevitable noises and the intense edges (*e.g.*, the bone boundary of trex) also cause disharmony. By training a high-resolution 2D PST network, UPST [7] is able to preserve the structure well. However, the color style of stylized scene is different from the reference. Instead, LipRF completely eliminates the limitations of CCPL and WCT<sup>2</sup>, and the stylized results have clear advantages in photorealism and color style.

**Tanks and Temples dataset.** Since UPST [7] does not support this dataset, Figure 5 depicts the comparison between LipRF and 2D PST methods. Similarly, 2D PST methods result in cross-view inconsistency like the sky color on the playground, and strong noises like the messy colors on car and ground. After transforming radiance field with Lips-

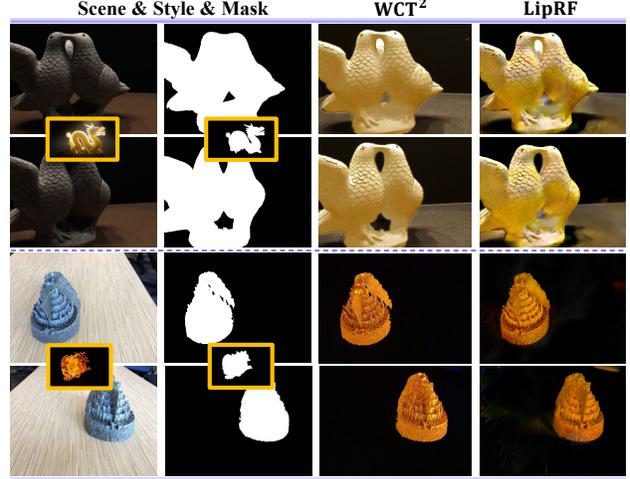


Figure 6. Object appearance editing. The semantic masks of style image and scene are provided to guide the style transfer.

chitz Network, our LipRF manages to alleviate these issues and obtains high-quality stylized results.

**Object appearance editing.** Another practical application of PST is object appearance editing [38], which leverages additional semantic masks of objects to guide style transfer. However, it is difficult to precisely annotate the object masks of each view. In Figure 6, the first scene comes from the DTU dataset [20], and the provided inaccurate masks result in the inconsistency of stylized images. The same problem is also observed for the second scene taken from LLFF [42] with automatic annotation. In contrast, LipRF obtains photorealistic and consistent results.

**Style interpolation.** Once the training of LipRF is completed, we can interpolate the source and stylized radiance fields to obtain  $\mathbf{F}_{sh}^\alpha = \alpha \mathbf{F}'_{sh} + (1 - \alpha) \mathbf{F}_{sh}$  by adjusting the factor  $\alpha \in [0, 1]$  (see Figure 7). This provides an efficient way to control or serialize the style change of scenes, saving much time compared with the interpolation of image pixels.

## 6.2. Quantitative results

**Consistency.** We take the temporal consistency [25] as the metric for evaluating the cross-view consistency. Specifically, given two view synthesis  $x_i$  and  $x_j$  of a scene, the temporal consistency is computed by

$$TC(x_i, x_j) = \frac{1}{|\mathcal{O}_{i,j}|} \|\mathcal{O}_{i,j} \mathcal{W}_{i,j}(x_i) - \mathcal{O}_{i,j} x_j\|^2, \quad (14)$$

where  $\mathcal{W}_{i,j}$  warps  $x_i$  to  $x_j$  according to the optical flow estimated by RAFT [56], and mask  $\mathcal{O}_{i,j}$  labels non-occluded pixels [53] in  $x_i$  and  $x_j$ .  $|\mathcal{O}_{i,j}|$  is the sum of tensor items. We further convert  $TC$  to the readable PSNR form

$$TC_{psnr}(x_i, x_j) = -10 * \log_{10}(TC(x_i, x_j)). \quad (15)$$

The metric is conducted on 8 scenes of LLFF dataset [42], and each scene will be stylized based on 4 images from



Figure 7. Style interpolation within radiance fields.

	fern	flower	fortress	horns	leaves	orchids	room	trex	Avg.
MKL	28.4	32.9	<b>34.2</b>	31.1	<b>26.2</b>	27.0	28.8	27.7	30.0
WCT <sup>2</sup>	24.0	28.5	24.1	27.9	22.7	24.3	25.9	24.7	25.3
w/ Lip	<b>29.5</b>	<b>34.1</b>	33.4	<b>32.6</b>	25.8	<b>28.0</b>	<b>29.4</b>	<b>28.7</b>	<b>30.3</b>
CCPL	22.0	23.4	22.5	25.1	20.5	21.1	23.9	24.7	22.9
w/ Lip	26.7	30.7	31.0	30.9	<b>26.2</b>	262	28.8	27.4	28.5
UPST	27.6	33.4	31.0	30.5	<b>26.2</b>	26.7	30.4	27.7	28.3

---

MKL	25.6	26.3	<b>28.0</b>	26.3	<b>22.9</b>	23.6	24.6	23.8	25.1
WCT <sup>2</sup>	22.1	23.7	19.1	24.4	20.0	21.5	22.2	21.4	21.8
w/ Lip	<b>26.7</b>	<b>28.0</b>	27.9	<b>27.8</b>	22.5	<b>24.6</b>	25.3	<b>24.3</b>	<b>25.9</b>
CCPL	20.5	19.8	18.3	22.5	18.7	18.9	21.5	22.4	20.3
w/ Lip	24.1	24.6	24.1	26.3	22.9	22.9	24.7	23.6	24.2
UPST	24.8	26.6	22.7	26.0	21.2	23.6	<b>25.7</b>	23.7	24.3

Table 1. Comparisons of short (upper) and long (lower) temporary consistency on LLFF dataset. The higher the value, the better the consistency. “w/ Lip” means coupling LipRF with the 2D PST method in the block. The last column shows the average value.

PST dataset [38]. The number of evaluated views is the same as that of training views. We report the short temporal consistency taking  $(x_i, x_{i+1})$  as the input pair, and the long temporal consistency taking  $(x_i, x_{i+5})$  as input pair. Table 1 details the average results and we have three main observations. The first is that LipRF is the only method that manifests similar or advanced properties over linear method MKL [51]. Second, LipRF stably promotes the consistency of 2D PST, thereby upgrading 2D PST to adapt for the 3D scene. Finally, LipRF does benefit from the improvement of 2D PST method, as evidenced by exhibiting better consistency of WCT<sup>2</sup> with LipRF against CCPL with LipRF.

**User study.** We further conduct subjective evaluation to compare the style effects. We operate LipRF with WCT<sup>2</sup>, UPST and MKL on 16 pairs of scenes and style images, then invite 35 volunteers to comprehensively assess and vote for their favorite rendered videos. Finally, LipRF obtains 79% preference that surpasses the 11% of UPST and 10% of MKL, proving the impressive effects of LipRF.

### 6.3. Ablation study

**Alternative activations.** The activation is essential for constructing Lipschitz networks, and the commonly used ReLU (or LeakyReLU), Sigmoid, Sine and Tanh are all Lipschitz continuous. In the experiments, we found that Sigmoid and Tanh do not work due to the vanishing gradients sometimes.

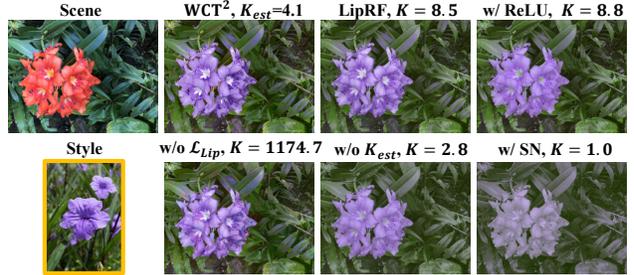


Figure 8. The results with different Lipschitz constants.  $K_{est}$  is the estimated value in Eq. (11). “w/ ReLU” replaces the sinusoidal activation of LipRF with ReLU. “w/o  $K_{est}$ ” sets the lower bound to zero. “w/ SN” replaces the adaptive regularization with the spectral normalization that forces the Lipschitz constant to 1.

In the remaining alternatives, the Sine-based LipRF can obtain more faithful color style compared with the ReLU-based one. Figure 8 visualizes the comparisons.

**Lipschitz regularization.** Next, we study the impact of the proposed adaptive Lipschitz regularization for LipRF. As shown in Figure 8, LipRF will create noisy details when removing the regularization. Meanwhile, the Lipschitz constant surges from 8.5 to 1174.7. Furthermore, when removing or setting  $K_{est}$  to zero, the Lipschitz constant of LipRF will decrease dramatically from 8.5 to 2.8 after 300 epochs. Both qualitative results validate our strategy. We also compare our strategy with spectral normalization [44], and the results indicate that the lower the Lipschitz constant, the weaker the color change, and the visually smoother the image. The comparison again proves the merits of LipRF.

## 7. Conclusion

In this paper, we present a novel and well-motivated framework namely LipRF to tackle photorealistic 3D scene stylization, where the main challenge is the preservation of photorealism and consistency during stylization. To this end, we first show that the two objectives can be simplified into maintaining volume rendering variance before and after stylization. Then, we prove that a Lipschitz MLP enables controlling the variance. Third, we propose an adaptive regularization to constrain the lower bound of Lipschitz constant and introduce the gradual gradient aggregation to optimize LipRF in a cost-efficient manner. Extensive experiments shows the versatility of LipRF. The main limitation is that, LipRF relies on the pre-trained radiance field. Once the radiance field cannot reconstruct the scene well, the learning of LipRF would be failed. We think this problem can be solved with the development of radiance fields.

**Acknowledgements.** This paper is supported by the National key research and development program of China (2021YFA1000403), and the National Natural Science Foundation of China (Nos. U19B2040, 11991022).

## References

- [1] Cem Anil, James Lucas, and Roger Grosse. Sorting out lipschitz function approximation. In *ICML*, 2019. 4
- [2] Soonmin Bae, Sylvain Paris, and Frédo Durand. Two-scale tone management for photographic look. In *ACM CGIT*, 2006. 2
- [3] Jonathan T. Barron. Squareplus: A softplus-like algebraic rectifier. *arXiv preprint arXiv:2112.11687*, 2021. 5
- [4] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. 2
- [5] Jiawen Chen, Andrew Adams, Neal Wadhwa, and Samuel W. Hasinoff. Bilateral guided upsampling. *ACM TOG*, 2016. 2, 3
- [6] Yang Chen, Yingwei Pan, Ting Yao, Xinmei Tian, and Tao Mei. Mocycle-gan: Unpaired video-to-video translation. In *ACM MM*, 2019. 2
- [7] Yaosen Chen, Qinjian Yuan, Zhiqiang Li, Yuegen Liu, Wen Wang, Chaoping Xie, Xuming Wen, and Qien Yu. Upst-nerf: Universal photorealistic style transfer of neural radiance fields for 3d scene. *arXiv preprint arXiv:2208.07059*, 2022. 6, 7
- [8] Pei-Ze Chiang, Meng-Shiun Tsai, Hung-Yu Tseng, Wei-Sheng Lai, and Wei-Chen Chiu. Stylizing 3d scene via implicit representation and hypernetwork. In *WACV*, 2022. 3, 5
- [9] Tai-Yin Chiu and Danna Gurari. Photowct2: Compact autoencoder for photorealistic style transfer resulting from blockwise training and skip connections of high-frequency residuals. In *WACV*, 2022. 2
- [10] Paul Debevec, Camillo J. Taylor, and Jitendra Malik. Modeling and rendering architecture from photographs: a hybrid geometry- and image-based approach. In *ACM CGIT*, 1996. 2
- [11] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR*, 2016. 1, 2, 3
- [12] Michaël Gharbi, Jiawen Chen, Jonathan T. Barron, Samuel W. Hasinoff, and Frédo Durand. Deep bilateral learning for real-time image enhancement. *ACM TOG*, 2017. 3
- [13] Georgia Gkioxari, Olivia Wiles, Richard Szeliski, Justin Johnson, Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *CVPR*, 2020. 2
- [14] Henry Gouk, Eibe Frank, Bernhard Pfahringer, and Michael J. Cree. Regularisation of neural networks by enforcing lipschitz continuity. *Machine Learning*, 2021. 3
- [15] Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of wasserstein gans. In *NeurIPS*, 2017. 3
- [16] Hsin-Ping Huang, Hung-Yu Tseng, Saurabh Saini, Maneesh Singh, and Ming-Hsuan Yang. Learning to stylize novel views. In *ICCV*, 2021. 3
- [17] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *CVPR*, 2017. 3
- [18] Yihua 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 *CVPR*, 2022. 3
- [19] Jing Huo, Shiyin Jin, Wenbin Li, Jing Wu, Yu-Kun Lai, Yinghuan Shi, and Yang Gao. Manifold alignment for semantically aligned style transfer. In *ICCV*, 2021. 1, 3
- [20] Rasmus Ramsbøl Jensen, A. Dahl, George Vogiatzis, Engil Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *CVPR*, 2014. 6, 7
- [21] Nikolai Kalischek, Jan Dirk Wegner, and Konrad Schindler. In the light of feature distributions: moment matching for neural style transfer. In *CVPR*, 2021. 2
- [22] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2015. 5
- [23] A. Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: benchmarking large-scale scene reconstruction. *ACM TOG*, 2017. 6, 7
- [24] Nicholas I. Kolkin, Jason Salavon, and Gregory Shakhnarovich. Style transfer by relaxed optimal transport and self-similarity. In *CVPR*, 2019. 2
- [25] Wei-Sheng Lai, Jia-Bin Huang, Oliver Wang, Eli Shechtman, Ersin Yumer, and Ming-Hsuan Yang. Learning blind video temporal consistency. In *ECCV*, 2018. 7
- [26] Anat Levin, Dani Lischinski, and Yair Weiss. A closed form solution to natural image matting. In *CVPR*, 2006. 2
- [27] Ming Li, Chunyang Ye, and Wei Li. High-resolution network for photorealistic style transfer. *arXiv preprint arXiv:1904.11617*, 2019. 2
- [28] Xueting Li, Sifei Liu, Jan Kautz, and Ming-Hsuan Yang. Learning linear transformations for fast image and video style transfer. In *CVPR*, 2019. 3
- [29] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. In *NeurIPS*, 2017. 1, 3
- [30] Yijun Li, Ming-Yu Liu, Xueting Li, Ming-Hsuan Yang, and Jan Kautz. A closed-form solution to photorealistic image stylization. In *ECCV*, 2018. 1, 2
- [31] Yanghao Li, Naiyan Wang, Jiaying Liu, and Xiaodi Hou. Demystifying neural style transfer. In *IJCAI*, 2017. 2
- [32] Chen-Hsuan Lin, Chen Kong, and Simon Lucey. Learning efficient point cloud generation for dense 3d object reconstruction. In *AAAI*, 2017. 2
- [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 3, 7
- [34] Hsueh-Ti Derek Liu, Francis Williams, Alec Jacobson, Sanja Fidler, and Or Litany. Learning smooth neural functions via lipschitz regularization. In *ACM SIGGRAPH*, 2022. 3
- [35] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas M. Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM TOG*, 2019. 2

- [36] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *ICLR*, 2017. 6
- [37] Ming Lu, Hao Zhao, Anbang Yao, Yurong Chen, Feng Xu, and Li Zhang. A closed-form solution to universal style transfer. In *ICCV*, 2019. 3
- [38] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. Deep photo style transfer. In *CVPR*, 2017. 1, 2, 4, 6, 7, 8
- [39] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021. 2
- [40] Nelson L. Max. Optical models for direct volume rendering. *IEEE TVCG*, 1995. 1, 2
- [41] Roey Mechrez, Eli Shechtman, and Lihi Zelnik-Manor. Photorealistic style transfer with screened poisson equation. In *BMVC*, 2017. 1, 2
- [42] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM TOG*, 2019. 2, 6, 7
- [43] 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. 1, 2, 3, 4, 6
- [44] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *ICLR*, 2018. 2, 3, 5, 8
- [45] Guido Montúfar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio. On the number of linear regions of deep neural networks. In *NeurIPS*, 2014. 4
- [46] Fangzhou Mu, Jian Wang, Yicheng Wu, and Yin Li. 3d photo stylization: Learning to generate stylized novel views from a single image. In *CVPR*, 2022. 3
- [47] Thu Nguyen-Phuoc, Feng Liu, and Lei Xiao. Snerf: Stylized neural implicit representations for 3d scenes. In *ACM SIGGRAPH*. 3, 5
- [48] Ido Omer and Michael Werman. Color lines: image specific color representation. In *CVPR*, 2004. 1, 2
- [49] Li Pan, Lei Zhao, Duanqing Xu, and Dongming Lu. Optimal transport of deep feature for image style transfer. In *ACM MM*, 2019. 2
- [50] François Pitié, Anil Kokaram, and Rozenn Dahyot. N-dimensional probability density function transfer and its application to color transfer. In *ICCV*, 2005. 2
- [51] François Pitié and Anil C. Kokaram. The linear monge-kantorovitch linear colour mapping for example-based colour transfer. In *CVMP*. 1, 4, 5, 6, 7, 8
- [52] Erik Reinhard, M. Adhikhmin, Bruce Gooch, and Peter Shirley. Color transfer between images. *IEEE CGA*, 2001. 1, 2, 4
- [53] Manuel Ruder, Alexey Dosovitskiy, and Thomas Brox. Artistic style transfer for videos. *Springer PR*, 2016. 7
- [54] Kevin Scaman and Aladin Virmaux. Lipschitz regularity of deep neural networks: analysis and efficient estimation. In *NeurIPS*, 2018. 3
- [55] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *NeurIPS*, 2020. 5
- [56] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020. 7
- [57] Kuan-Wei Tseng, Jing-Yuan Huang, Yang-Shen Chen, Chu-Song Chen, and Yi-Ping Hung. Pseudo-3d scene modeling for virtual reality using stylized novel view synthesis. In *ACM SIGGRAPH*, 2022. 1
- [58] Michael Waechter, Nils Moehrl, and Michael Goesele. Let there be color! large-scale texturing of 3d reconstructions. In *ECCV*, 2014. 2
- [59] Zijie Wu, Zhen Zhu, Junping Du, and Xiang Bai. Ccpl: Contrastive coherence preserving loss for versatile style transfer. In *ECCV*, 2022. 1, 3, 6, 7
- [60] Xide Xia, Tianfan Xue, Wei-Sheng Lai, Zheng Sun, Abby Chang, Brian Kulis, and Jiawen Chen. Real-time localized photorealistic video style transfer. In *WACV*, 2021. 3
- [61] Xide Xia, Meng Zhang, Tianfan Xue, Zheng Sun, Hui Fang, Brian Kulis, and Jiawen Chen. Joint bilateral learning for real-time universal photorealistic style transfer. In *ECCV*, 2020. 1, 3
- [62] Xuezhong Xiao and Lizhuang Ma. Color transfer in correlated color space. In *ACM VRCIA*, 2006. 4
- [63] Yao-Yuan Yang, Cyrus Rashtchian, Hongyang Zhang, Ruslan Salakhutdinov, and Kamalika Chaudhuri. A closer look at accuracy vs. robustness. In *NeurIPS*, 2020. 3
- [64] Kangxue Yin, Jun Gao, Maria Shugrina, Sameh Khamis, and Sanja Fidler. 3dstylenet: Creating 3d shapes with geometric and texture style variations. In *ICCV*, 2021. 3
- [65] Jaejun Yoo, Youngjung Uh, Sanghyuk Chun, Byeongkyu Kang, and Jung-Woo Ha. Photorealistic style transfer via wavelet transforms. In *ICCV*, 2019. 1, 2, 6
- [66] Yuichi Yoshida and Takeru Miyato. Spectral norm regularization for improving the generalizability of deep learning. *arXiv preprint arXiv:1705.10941*, 2017. 3
- [67] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*. 1, 2, 3, 4, 5
- [68] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 3, 4
- [69] Kai Zhang, Nicholas I. Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snavely. Arf: Artistic radiance fields. In *ECCV*, 2022. 3, 5
- [70] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2
- [71] Yabin Zhang, Minghan Li, Ruihuang Li, Kui Jia, and Lei Zhang. Exact feature distribution matching for arbitrary style transfer and domain generalization. In *CVPR*, 2022. 2
- [72] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: learning view synthesis using multiplane images. *ACM TOG*, 2018. 2