

---

# MM-NeRF: Multimodal-Guided 3D Multi-Style Transfer of Neural Radiance Field

---

**Zijiang Yang**

School of Automation and Electrical Engineering  
University of Science and Technology Beijing  
zijiangyang@xs.ustb.edu.cn

**Zhongwei Qiu**

School of Automation and Electrical Engineering  
University of Science and Technology Beijing  
qiuzhongwei@xs.ustb.edu.cn

**Chang Xu**

School of Computer Science  
University of Sydney  
c.xu@sydney.edu.au

**Dongmei Fu**

School of Automation and Electrical Engineering  
University of Science and Technology Beijing  
fdm\_ustb@ustb.edu.cn

## Abstract

3D style transfer aims to render stylized novel views of 3D scenes with the specified style, which requires high-quality rendering and keeping multi-view consistency. Benefiting from the ability of 3D representation from Neural Radiance Field (NeRF), existing methods mainly learn the stylized NeRF by giving a reference style from an image. However, they still suffer the challenges of high-quality stylization with texture details for multi-style transfer and stylization with multimodal guidance. In this paper, we reveal that the same objects in 3D scenes show various states (color tone, details, etc.) from different views after stylization since previous methods optimized by single-view image-based style loss functions, leading NeRF to tend to smooth the texture details, further resulting in low-quality rendering for 3D multi-style transfer. To tackle these problems, we propose a novel **Multimodal-guided 3D Multi-style transfer of NeRF**, termed **MM-NeRF**, which achieves high-quality 3D multi-style rendering with texture details and can be driven by multimodal-style guidance. First, MM-NeRF adopts a unified framework to project multimodal guidance into CLIP space and extracts multimodal style features to guide the multi-style stylization. To relieve the problem of lacking details, we propose a novel Multi-Head Learning Scheme (MLS) for multi-style transfer, in which each style head predicts the parameters of the color head of NeRF. MLS decomposes the learning difficulty caused by the inconsistency of multi-style transfer and improves the quality of stylization. In addition, the MLS can generalize pre-trained MM-NeRF to any new styles by adding heads with small training costs (a few minutes). Extensive experiments on three real-world 3D scene datasets show that MM-NeRF achieves high-quality 3D multi-style stylization with multimodal guidance, keeps multi-view consistency, and keeps semantic consistency of multimodal style guidance. Codes will be released later.

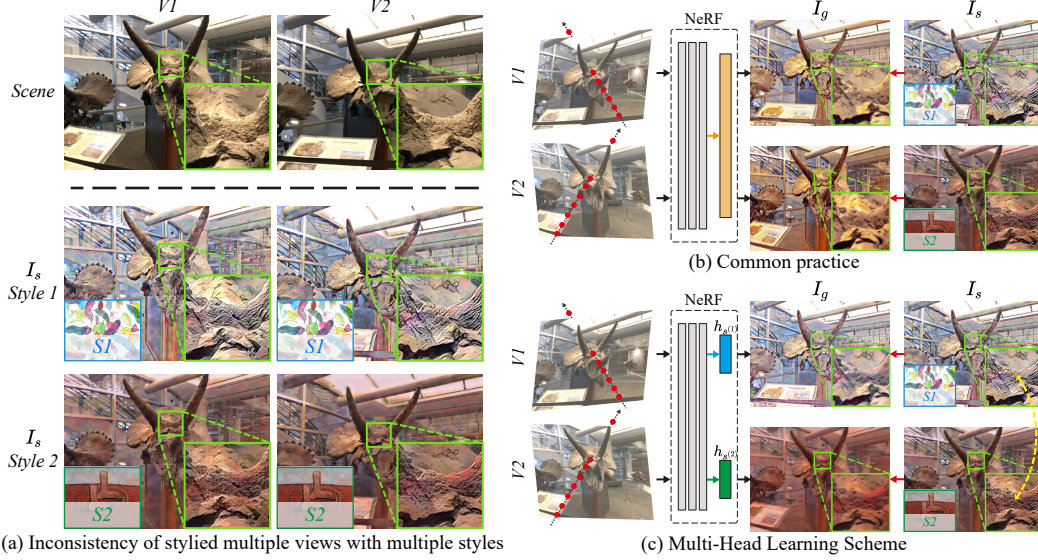


Figure 1: (a) Optimizing stylized NeRF by image-based style loss involves the challenges of i) inconsistent stylized images of multiple views with the same style and ii) inconsistent stylized images of the same view with multiple styles. (b) The common practice [67, 26, 34, 40, 13] with a single prediction head generates images with blurry texture details. (c) To reduce the optimization difficulty, MM-NeRF adopts a Multi-Head Learning Scheme to decompose the learning difficulty and improve the quality of texture details in generated images.  $V1$  and  $V2$  are two different views.  $S1$  and  $S2$  are two different styles.  $I_g$  and  $I_s$  are generated images by NeRF and stylized images by a 2D style transfer method, respectively.  $h_{s(1)}$  and  $h_{s(2)}$  are color heads of style  $s^{(1)}$  and style  $s^{(2)}$ , respectively.

## 1 Introduction

Recently, Neural Radiance Field [36] (NeRF) has been widely used in 3D representation since it can model complex 3D scenes with small costs. Thus, the NeRF-based downstream tasks, such as relighting [70, 53], surface reconstruction [59] and 3D object generation [30], have been rapidly developed, including 3D style transfer [67, 34, 26]. Given a reference style and multi-view images of a 3D scene, 3D style transfer aims to synthesize novel views of the scene that match the color tone and stroke of this style. Existing methods [11, 67, 40, 13, 26, 34] have made great progress on 3D style transfer tasks, such as photo-realistic style transfer [11], arbitrary style transfer [34, 13], and artistic style transfer [67, 40, 26]. However, these methods only focus on the stylization guided by image-based styles and still suffer the challenges of high-quality stylization with texture details for multi-style transfer. In this work, we aim to achieve high-quality 3D multi-style transfer of NeRF with texture details and multi-style stylization guided by multimodal styles (image, text, audio, etc.).

The style transfer of 3D scenes requires multi-view stylized images or features as the supervision to train the stylized NeRF model. Existing methods [40, 11, 26] usually transfer the multi-view images to stylized images by 2D stylization models [25, 54], which is further used as the supervision of 3D style transfer. Specifically, SNeRF [40] achieves 3D single-style transfer by giving an image style. While UPST-NeRF [11] and StylizedNeRF [26] utilize a set of pre-defined image-based styles to realize 3D multi-style transfer. Besides, StyleRF [34] fuses the image-based style features and 3D scene features to achieve zero-shot 3D style transfer. However, these methods only adopt single-modality styles (image-based styles) as guidance, which limits their application. In this work, we study using multimodal styles as the guidance of 3D style transfer.

In addition, existing methods [67, 11, 26] of 3D style transfer still remain the challenge of generating stylized images with texture details. As shown in Figure 1 (b), the stylized images  $I_g$  from existing practice show blurry details and limited matching style with the reference styles. The reason is that these methods are mainly optimized by single-view image-based style loss function and it is computed with view-inconsistent 2D style transfer methods. It means that the same objects in the supervision  $I_s$  show various states (color tone, details, etc.) from different views, and the inconsistency of optimization target results in the stylized NeRF model tend to smooth the details, further generating



low-quality stylization results. A more detailed analysis is in Section 4.1. In this work, we argue the inconsistency of optimization targets results in low-quality stylization and aim to relieve this problem.

To tackle the above two problems, we propose a framework of multimodal-guided 3D multi-style transfer of NeRF, termed MM-NeRF. In MM-NeRF, we break down the multimodal-guided 3D multi-style transfer into two components, including the feature extraction of multimodal styles and the learning of multiple styles. For feature extraction of multimodal styles, MM-NeRF transfers multimodal styles into images with pre-trained cross-modal generation and encodes these styles into a unified CLIP [48] feature space, which is used to guide the stylization further. For the learning of multiple styles, as the inconsistency of multi-view and multi-style optimization targets leads models to tend to smooth the generated details in images, MM-NeRF decomposes this mutually exclusive optimization problem into multiple sub-optimization problems and optimizes an independent prediction head for each sub-optimization problem by Multi-Head Learning Scheme (MLS). Therefore, the multimodal style extractor and Multi-Head Learning Scheme lead MM-NeRF to achieve high-quality stylization with texture details guided by multimodal styles. In addition, benefiting from the MLS, MM-NeRF can be generalized to any new style by adding new heads with small training costs (a few minutes since just need to train the new heads).

To the best of our knowledge, MM-NeRF is the first framework that achieves the multimodal-guided 3D multi-style transfer of NeRF. Our contributions can be summarized as follows:

- We propose MM-NeRF, the first unified framework of multimodal-guided 3D multi-style transfer of NeRF, which achieves high-quality 3D multi-style stylization with multimodal guidance (image, text, audio, etc.).
- We analyze the lack of details in existing 3D stylization schemes and reveal that it’s due to the inconsistency of optimization targets of 3D multi-style transfer. To tackle this problem, we propose a novel Multi-Head Learning Scheme (MLS) to decompose the optimization difficulty of inconsistency.
- We propose a new incremental learning mechanism based on MLS, which generalizes pre-trained MM-NeRF to any new styles with small training costs.
- Extensive experiments conducted on three widely-used datasets demonstrate that MM-NeRF outperforms prior methods and exhibits advantages in 3D consistency, multimodal guidance, and semantic consistency of multimodal guidance.

## 2 Related Work

### 2.1 Neural Radiance Field

Given multi-view images of a scene, NeRF [36] and its variants [66, 43, 5, 12, 21, 39] can render high quality novel views. Due to the ability to learn high-resolution complex scenes, NeRF receives widespread attention and has been extended to dynamic scenes, large-scale scenes [35, 69], relighting [70, 53, 51], surface reconstruction [59], generation tasks [41, 57, 8, 6], etc. In addition, in the field of dynamic scenes, for example 3D human pose[45, 46, 47], NeRF has made great progress[18, 63, 27]. In this paper, we focus on 3D style transfer, and NeRFs are employed to learn implicit representations of scenes.

### 2.2 Neural Style Transfer

**2D Style Transfer.** Neural image style transfer aims to transfer the style of a reference image to other images using neural networks [19]. For fast image style transfer, methods based on feed-forward networks formulate style transfer as an optimization problem of neural networks and have made significant progress [28, 25, 32, 15, 54, 4]. In recent years, neural style transfer has gradually expanded from images to videos, and stylized videos are generated with a reference style image [23, 50, 60, 10, 16].

**3D Style Transfer.** Given multi-view images of a scene and a reference style, 3D style transfer aims to synthetic stylized novel views of this scene with the specified style. Early 3D style transfer methods are mostly based on point clouds [7, 24, 37]. The resolution of 3D scenes limits the stylization quality of these methods, and this problem is particularly evident in real-world scenes. Recently, NeRF-based

3D style transfer develops rapidly [11, 67, 40, 13, 26, 34]. Zhang et al. [67] and Nguyen-Phuoc et al. [40] propose high-quality single-style transfer models, but time-consuming optimization is required for each style. Chiang et al. [13] is the first to formulate the multi-style transfer of NeRF as a weights prediction problem, but their method produces blurry results. Liu et al. [34] achieve zero-shot 3D style transfer of arbitrary style by fusing high-level features, but the stylization quality of new style can not be guaranteed, and stylized images lose texture details.

## 2.3 Multimodal Learning

Multimodal representation learning [42, 3, 20, 57] expands the application scope of generative models [33, 49, 52, 9, 22, 71]. By fusing complementary information from multimodal, the performance of related tasks also can be significantly improved. Li et al. [33] generates videos with text descriptions. Rombach et al. [49] and Saharia et al. [52] achieve high-quality image generation based on Diffusion Model with text. Chen et al. [9] and Hao et al. [22] propose cross-modal audio-visual generations, which can achieve four generation modes: audio-to-visual, visual-to-audio, audio-to-audio, and visual-to-visual. In the field of 3D generation, current multimodal works mainly focus on text-guided 3D object generation [57, 62, 58, 30, 44] and 3D style transfer with multimodal guidance has not been studied yet. In this paper, we tackle the multimodal-guided 3D multi-style transfer and propose the first unified framework to achieve it.

## 3 Preliminaries

NeRF [36] is proposed to model a scene as an implicitly volumetric field containing a density field representing objects and a radiation field representing texture. Given 3D position  $\mathbf{x} \in \mathcal{R}^3$  and viewing direction  $\mathbf{d} \in \mathcal{R}^2$ , NeRF uses Multi-Layer Perceptions (MLPs) to predict opacity  $\sigma(\mathbf{x}) \in \mathcal{R}^+$  and view-dependent radiance color  $\mathbf{c}(\mathbf{x}, \mathbf{d}) \in \mathcal{R}^3$ . Ray marching is used to visualize the implicitly volumetric field and the color of a pixel  $C(\mathbf{r})$  is determined by the integral:

$$C(\mathbf{r}) = \int_{t=0}^{\infty} \sigma(\mathbf{o} + t\mathbf{d}) \cdot \mathbf{c}(\mathbf{o} + t\mathbf{d}, \mathbf{d}) \cdot e^{-\int_{t=0}^t \sigma(\mathbf{o} + l\mathbf{d}) dl} dt, \quad (1)$$

where,  $\mathbf{r}$  is the ray passing through the pixel and  $\mathbf{o} \in \mathcal{R}^3$  is the 3D position of the camera. NeRF are optimized by minimizing losses between predicted colors according to Equation 1 and ground truth colors. After training, given a new viewpoint, NeRF synthesizes novel views of the same scene.

## 4 Method

### 4.1 Problem Formulation

**Multimodal-Guided 3D Multi-Style Transfer.** Given multi-view images of a scene with camera intrinsic and extrinsic parameters, the goal of multimodal-guided 3D multi-style transfer is to synthesize novel stylized views of scenes with multimodal styles. In subsequent discussions, we focus on the stylization of NeRF, and this task is as follows:

$$C(\mathbf{r}; s^{(i)}) = \int_{t=0}^{\infty} \sigma_s(\mathbf{o} + t\mathbf{d}; s^{(i)}) \cdot \mathbf{c}_s(\mathbf{o} + t\mathbf{d}, \mathbf{d}; s^{(i)}) \cdot e^{-\int_{t=0}^t \sigma_s(\mathbf{o} + l\mathbf{d}; s^{(i)}) dl} dt, s^{(i)} \in S, \quad (2)$$

where  $S$  is the style set,  $N_s$  is the number of styles in  $S$ ,  $s^{(i)}$  is the  $i$ th style in style set  $S = \{s^{(i)}\}_{i=1}^{N_s}$ ,  $\sigma_s$  is the stylized opacity field and  $\mathbf{c}_s$  is the stylized radiation field. For multi-style,  $N_s > 1$  and, for multimodal,  $s^{(i)}$  can be an image-based style, a text-based style, etc. the main task of Multimodal-Guided 3D Multi-Style Transfer of NeRF is to determine and optimize  $\sigma_s$  and  $\mathbf{c}_s$ .

**The Problem of Inconsistency.** Existing 3D style transfer methods are mainly optimized with single-view image-based style loss functions. Especially, stylized images and features of views, which are used to supervise the learning of 3D stylization models, are generated by 2D style transfer methods. However, 2D style transfer methods do not consider consistency across multiple views, which leads to two inconsistency issues, including (1) inconsistent stylized images of multiple views generated from 2D stylization methods with the same style and (2) inconsistent stylized images of the same view generated from 2D stylization methods with multiple styles. As illustrated in Figure

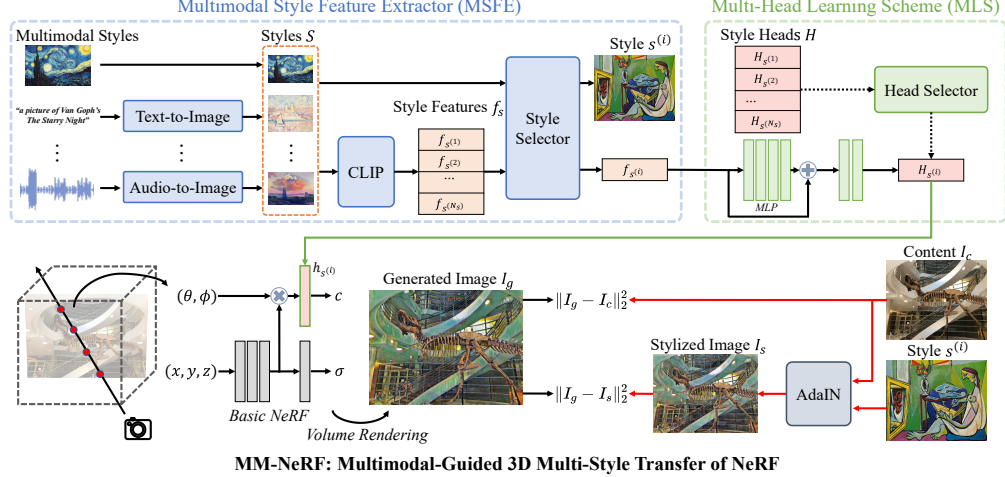


Figure 2: The framework of MM-NeRF includes a basic NeRF, MSFE, and MLS. MSFE transfers multimodal styles into images and creates the set of style images  $S$ . Then, MSFE selects a style  $s^{(i)}$  and encodes the style into CLIP [48] space to generate style features  $f_s^{(i)}$ . MLS maintains a prediction head for each style and predicts the parameters of the color head in NeRF from the style feature. By replacing the color head in basic NeRF with the predicted parameters, MM-NeRF achieves high-quality 3D multi-style transfer with multimodal guidance.  $\otimes$ ,  $\oplus$ , AdaIN [25] are sum, concatenation, and 2D stylization method, respectively.

1 (a), color tone and texture details of the same object are view-dependent and style-dependent. In 3D multi-style transfer, these two types of inconsistency together affect parameter optimization, increasing optimization difficulty. The optimization problem of multi-style 3D style transfer of NeRF is as follows:

$$\{\sigma_s, c_s\} = \arg \min_{\{\sigma_s, c_s\}} \sum_{j=1}^{N_s} \sum_{i=1}^{N_{\mathcal{R}}} \|C(\mathbf{r}^{(i)}; s^{(j)}) - C_{s^{(j)}}(\mathbf{r}^{(i)})\|_2^2, \quad (3)$$

where  $\mathcal{R}$  is the set of rays,  $N_{\mathcal{R}}$  is the total number of rays,  $\mathbf{r}^{(i)}$  is the  $i$ th ray in  $\mathcal{R}$ ,  $C_s(\mathbf{r}^{(i)})$  is the 2D stylized color of  $\mathbf{r}^{(i)}$ . Although NeRF can forcibly ensure consistency across multi-view, due to the high inconsistency of  $C_s(\mathbf{r}^{(i)})$ , as shown in Figure 1 (b), NeRF tends to smooth texture details to reduce the total error, resulting in low-quality stylized novel views.

**Multi-Head Learning Scheme.** The inconsistencies mentioned above can be summarized as the inconsistency of optimization targets, constituting a mutually exclusive multi-objective optimization problem. To address this problem, we propose Multi-Head Learning Scheme (MLS) to decompose the mutually exclusive multi-objective optimization problem into multiple sub-optimization problems and optimize an independent sub-module for each sub-optimization problem.

Specifically, in this paper, the prediction model of radiance color is decomposed into two parts, including a fully connected network shared by all styles and a prediction head for each style:

$$c_s(\mathbf{o} + t\mathbf{d}, \mathbf{d}; s^{(i)}) = h_{s^{(i)}}(b_s(\mathbf{o} + t\mathbf{d}, \mathbf{d})), s^{(i)} \in S, \quad (4)$$

where  $h_{s^{(i)}}$  is the color head of style  $s^{(i)}$ ,  $b_s$  is the shared fully connected network of all styles in stylized radiation field  $c_s$ . Since  $h_{s^{(i)}}$  only suffers from the inconsistency caused by multi-view, the model can be optimized to a better solution. As shown in Figure 1 (c), compared to NeRF with a single head in Figure 1 (b), MLS can significantly enhance the clarity of texture details of generated images.

## 4.2 MM-NeRF

The framework of MM-NeRF is illustrated in Figure 2, which includes a basic NeRF, Multimodal Style Feature Extractor (MSFE), and Multi-Head Learning Scheme (MLS).

**Basic NeRF.** Similar to the original NeRF [36], in MM-NeRF, the stylized opacity field  $\sigma_s$  and the stylized radiation field  $c_s$  share a backbone  $F_b$ . The output of  $F_b$  is mapped to a vector by a single

---

**Algorithm 1** Stylization Training of MM-NeRF

---

**Input:** Style set  $S$ , multi-view images of a real-world scene  $\{I_c^{(i)}\}_{i=1}^{N_c}$ , NeRF pre-trained on  $\{I_c^{(i)}\}_{i=1}^{N_c}$ , Multimodal Style Feature Extractor (MSFE), the Multi-Head Learning Scheme (MLS).

**Parameter:** Iteration steps of Stylization training  $T_s$ , batch size of each step  $M$ .

**Stylization Training**

- 1: **for** iteration  $t = 1, \dots, T_s$  **do**
  - 2:   Randomly sample an image  $I_c$  from  $\{I_c^{(i)}\}_{i=1}^{N_c}$  and a style  $s$  from  $S$ .
  - 3:   Extract the feature vector  $f_s$  by MSFE according to Equation 5.
  - 4:   Predict parameters  $w_s$  of the color head of NeRF by MLS according to Equation 6.
  - 5:   Randomly sample a mini-batch  $\{\mathbf{r}^{(i)}\}_{i=1}^M$  of rays from  $I_c$ .
  - 6:   Predict colors  $\hat{C}(\mathbf{r}^{(i)}; s)$ ,  $\mathbf{r}^{(i)} \in \{\mathbf{r}^{(i)}\}_{i=1}^M$  according to Equation 2.
  - 7:   Optimize MLS to minimize stylization training loss  $\mathcal{L}$  according to Equation 7.
  - 8:   Update the parameters of MLS.
  - 9: **end for**
- 

full connection layer. An exponential function activates the first element of this vector to predict opacity, and this vector is also used as input of the color head  $F_c$  to predict RGB radiance.

**Multimodal Style Feature Extractor.** MSFE transfers multimodal styles into images to capture the color tone, stroke, and texture details of styles. For any style  $s^{(i)}$  in the style set  $S$ , the feature extraction is as following:

$$f_{s^{(i)}} = E(T(s^{(i)})), s^{(i)} \in S, \quad (5)$$

where  $f_{s^{(i)}}$  is the feature vector of  $s^{(i)}$ ,  $T$  is the modal transformation model and  $E$  is the encoder. Stable Diffusion [49] is employed for text-to-image transformation, and image-based styles do not require any transformation operations. In addition, Contrastive Language-Image Pre-Training (CLIP) [48] is employed as an encoder to extract style features. To verify the extensibility of MSFE, we also pre-train an audio-to-image model to achieve audio-guided 3D style transfer.

**Learning of Multi-Style.** MM-NeRF learns multiple styles by MLS. We implement a parameter predictor to predict parameters of the color head of NeRF with multiple styles:

$$w_{s^{(i)}} = H_{s^{(i)}}(B(f_{s^{(i)}})), s^{(i)} \in S, \quad (6)$$

where  $w_{s^{(i)}}$  is the predicted parameters of the color head  $h_{s^{(i)}}$ ,  $H_{s^{(i)}}$  is the prediction head of style  $s^{(i)}$ ,  $B$  is a fully connected network shared by all styles in MLS.

### 4.3 Learning

The training process is shown as Algorithm 1. We first pre-train NeRF using multi-view images for each scene and fix it in the stylization training. Multi-view images are denoted as  $\{I_c^{(i)}\}_{i=1}^{N_c}$ , where  $I_c^{(i)}$  is the  $i$ th image and  $N_c$  is the number of images of this scene. At each iteration of stylization training, we randomly sample an image  $I_c$  from  $\{I_c^{(i)}\}_{i=1}^{N_c}$  and a style from the set of styles  $S$ . Parameters of the color head of NeRF are predicted as described in Section 4.2 and NeRF predicts the stylized color of pixels. MLS is optimized to minimize stylization training loss.

**Loss Functions.** To reduce the computational complexity of each step, we randomly sample  $M$  pixels from the supervised image and a mini-batch of rays  $\{\mathbf{r}^{(i)}\}_{i=1}^M$  of these pixels is created to render predicted colors with stylized NeRF. The stylization training loss function  $\mathcal{L}$  is as follows:

$$\mathcal{L} = \mathcal{L}_s + \lambda_c \mathcal{L}_c, \quad (7)$$
$$s.t. \mathcal{L}_s = \frac{1}{M} \sum_{i=1}^M \|\hat{C}(\mathbf{r}^{(i)}; s) - I_s(\mathbf{r}^{(i)})\|_2^2, \mathcal{L}_c = \frac{1}{M} \sum_{i=1}^M \|\hat{C}(\mathbf{r}^{(i)}; s) - I_c(\mathbf{r}^{(i)})\|_2^2,$$

where  $\mathcal{L}_s$  is the style loss,  $\mathcal{L}_c$  is the content loss,  $\lambda_c$  is the loss weight of the content loss,  $\hat{C}(\mathbf{r}^{(i)}; s)$  is the predicted color of ray  $\mathbf{r}^{(i)} \in \{\mathbf{r}^{(i)}\}_{i=1}^M$  and  $I_s$  is the stylized image of sampled photo-realistic image  $I_c$ . Pre-trained AdaIN [25] is based on feature fusion to support arbitrary style transfer and is employed to generate the stylized image  $I_s$  based on style  $T(s)$  and content image  $I_c$ :

$$I_s = \text{AdaIN}(T(s), I_c). \quad (8)$$

#### 4.4 Incremental Learning Scheme of New Styles

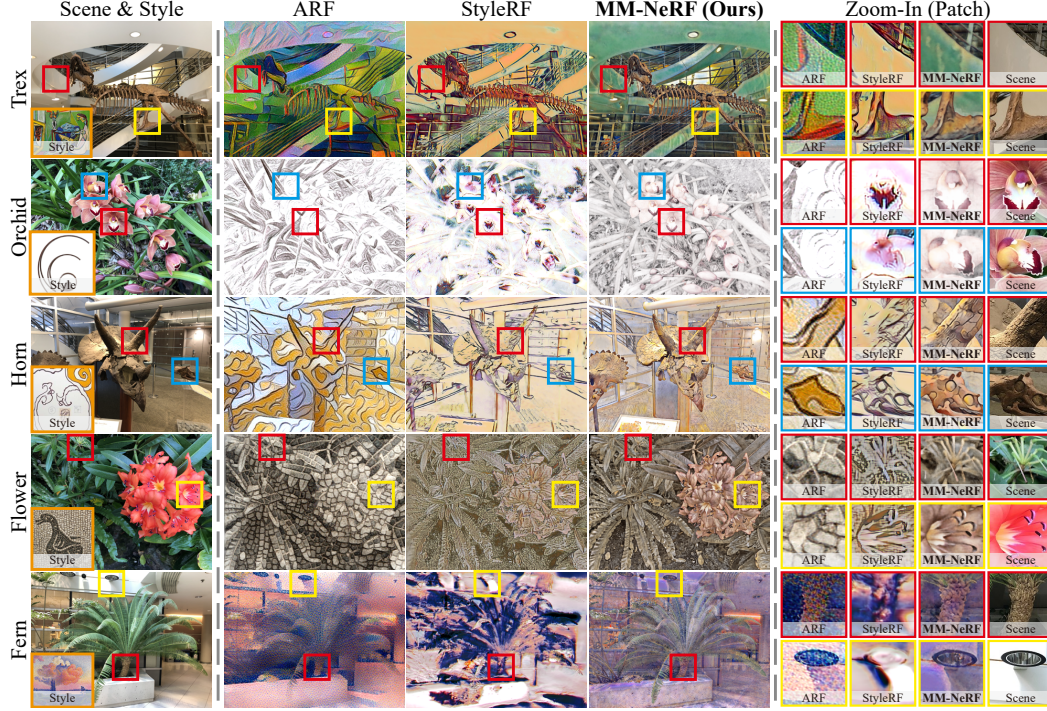


Figure 3: Comparison of image-guided 3D style transfer with ARF[67] and StyleRF[34] on forward-facing scenes[36]. MM-NeRF generates high-quality stylized images with texture details.

Method	Fern	Flower	Horn	Orchid	Trex	Average(↑)
ARF [67]	0.3119	0.2058	0.2203	0.2283	0.2583	0.2449
StyleRF [34]	0.2712	0.3217	0.3468	0.1748	0.4424	0.3114
<b>MM-NeRF (Ours)</b>	<b>0.5377</b>	<b>0.4682</b>	<b>0.5626</b>	<b>0.3910</b>	<b>0.5990</b>	<b>0.5117</b>

Table 1: Comparison of stylized details with SOTA methods on SSIM(↑).

Method	Short-range Consistency		Long-range Consistency	
	TWE(↓)	LPIPS(↓)	TWE(↓)	LPIPS(↓)
ARF [67]	0.0090	0.0694	0.0231	0.2694
StyleRF [34]	0.0072	0.0599	0.0121	0.2407
<b>MM-NeRF (Ours)</b>	<b>0.0025</b>	<b>0.0293</b>	<b>0.0060</b>	<b>0.1284</b>

Table 2: Comparison of consistency with SOTA methods on warped LPIPS(↓) and TWE(↓).

**Similarity Measurement.** For a new style  $s^{(new)}$ , we first measure similarity between  $s^{(new)}$  and other styles in the style set  $S$  on CLIP space by cosine distance:

$$\mathbf{f}_{s^{(sim)}} = \min(1 - \frac{\mathbf{f}_{s^{(i)}}^T \cdot \mathbf{f}_{s^{(new)}}}{\|\mathbf{f}_{s^{(i)}}\|_2 \cdot \|\mathbf{f}_{s^{(new)}}\|_2}), s^{(i)} \in S, \quad (9)$$

where  $\mathbf{f}_{s^{(new)}}$  and  $\mathbf{f}_{s^{(i)}}$  are the feature vectors of  $s^{(new)}$  and  $s^{(i)}$ , respectively.  $s^{(sim)}$  is the most similar style to the new style in the style set and  $\mathbf{f}_{s^{(sim)}}$  is the feature vector of  $s^{(sim)}$ . If the cosine distance between  $s^{(sim)}$  and  $s^{(new)}$  is smaller than the threshold  $\gamma$ , it is determined that the style is already included in the style set and the prediction head of  $s^{(sim)}$  in MLS will be called for stylization. Otherwise, incremental learning for style  $s^{(new)}$  will be performed.

**Incremental Learning.** For a new style, a new prediction head is initialized with the prediction head of  $s^{(sim)}$ , which is determined by similarity measurement. The incremental learning of new styles is similar to the stylization training described in Section 4.3. The only difference is that the network shared by all styles in MLS is fixed during incremental learning. In particular, as only new prediction heads need to be optimized, incremental learning of new styles can be finished in several minutes.



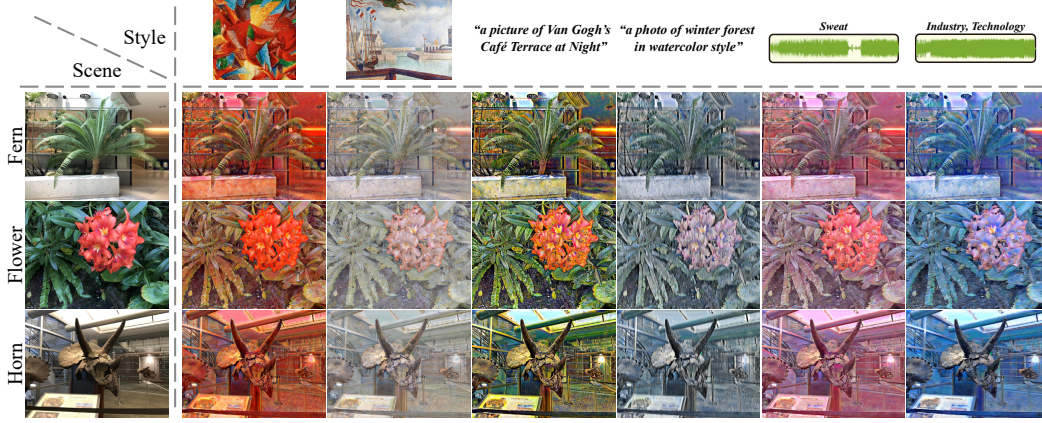


Figure 4: Multimodal-Guided 3D Style Transfer. MM-NeRF achieves multimodal-guided 3D style transfer on forward-facing scenes [36].

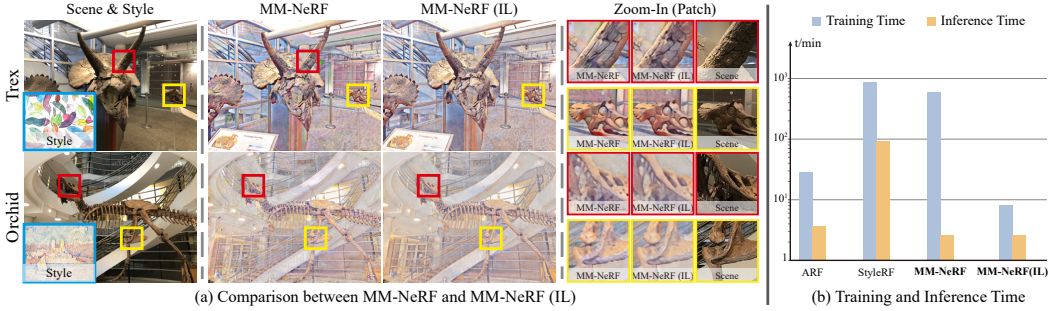


Figure 5: Incremental Learning of MM-NeRF. (a) Incremental Learning on Trex and Orchid [36]. (b) Comparison of training time and inference time with ARF[67] and StyleRF[34].

Method	Consistency	Stylization		
		Overall	Color Tone	Texture Details
ARF [67]	17.47%	33.14%	40.94%	4.67%
StyleRF [34]	12.16%	13.42%	12.80%	22.56%
<b>MM-NeRF (Ours)</b>	<b>70.37%</b>	<b>53.44%</b>	<b>46.26%</b>	<b>72.77%</b>

Table 3: The user study compares MM-NeRF with SOTA methods in consistency and stylization.

## 5 Experiments

**Datasets & Metrics.** We conduct extensive experiments on three widely-used 3D scene datasets [36, 29, 65] to evaluate MM-NeRF. Images from WikiArt [2], self-written texts, and audio from Metal Albums Artwork [1] are used to guide style transfer.

3D style transfer models are evaluated on multi-view consistency and stylized quality. In our experiments, we follow the usual practice to evaluate short-range consistency and long-range consistency by Temporal Warping Error (TWE) and warped Learned Perceptual Image Patch Similarity (LPIPS) [24, 13, 40]. Stylized quality is mainly evaluated by a user study. Structural Similarity (SSIM) between stylized views and photo-realistic views is also employed to evaluate details.

**Baselines.** We compare MM-NeRF to state-of-the-art NeRF-based stylization methods [67, 34]. ARF [67] is a typical single-style stylization model and StyleRF [34] supports style transfer of arbitrary image style. As MM-NeRF is the first multimodal-guided 3D style transfer model, we mainly compare MM-NeRF with these methods on image-guided 3D style transfer.

**Experiment Settings.** MM-NeRF is trained on pre-defined styles from WikiArt [2]. Text-guided 3D style transfer and audio-guided 3D style transfer are implemented by incremental learning described in Section 4.4. For each scene, NeRF is pre-trained for 20k iterations and MM-NeRF is trained for 60k iterations on a single NVIDIA Tesla P40 GPU. The hyper-parameters of  $\lambda_c$  and  $\gamma$  are set to 1e-3 and 0.95, respectively. ARF [67] and StyleRF [34] are implemented by their released code.



Figure 6: Comparison of semantic consistency of multimodal styles on Ship [65] and Truck [28].

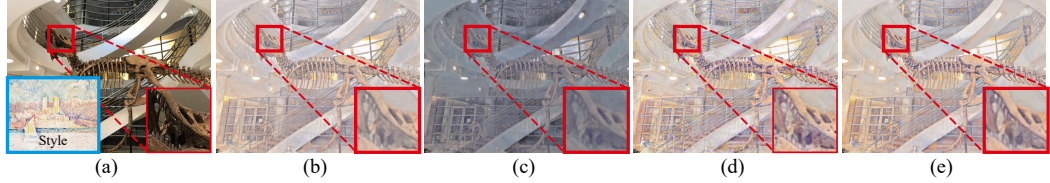


Figure 7: Ablation Study of MLS in MM-NeRF. (a) Scene and Style. (b) The stylized novel view of pre-defined style by single-head parameter predictor. (c) The stylized novel view of new style by single-head parameter predictor. (d) The stylized novel view of pre-defined style by MM-NeRF. (e) The stylized novel view of new style by MM-NeRF (Incremental Learning).

## 5.1 Main Results

**High-quality 3D Multi-style Transfer with Details.** In Figure 3, we compare MM-NeRF with ARF [67] and StyleRF [34] on forward-facing scenes [36]. MM-NeRF generates high-quality stylized novel views with texture details. Although StyleRF [34] and ARF [67] can generate highly stylized images, severe loss of details makes it difficult to recognize objects. Table 1 reports qualitative results of SSIM. Compared to other methods, MM-NeRF preserves more structural details. Table 2 reports the comparisons on consistency and MM-NeRF significantly outperforms other methods.

**Multimodal-Guided 3D Style Transfer.** Figure 4 reports multimodal-guided 3D style transfer of MM-NeRF on forward-facing scenes [36]. MM-NeRF can generate stylized novel views with reasonable color tones and details with the guidance of images and text. In addition, MM-NeRF can achieve audio-guided style 3D transfer by introducing a pre-trained audio-to-image model.

**Incremental Learning of New Styles.** MM-NeRF can generalize to new styles by Incremental Learning (IL). To evaluate the quality of new styles learned by IL, we first remove a style image from the style set. MM-NeRF is trained with the remaining styles and the specified style image is learned by IL. As shown in Figure 5 (a), MM-NeRF can generate similarly stylized novel views with pre-defined styles and new styles. Figure 5 (b) reports the comparison of training time and inference time with ARF [67] and StyleRF [34]. Although MM-NeRF requires large training costs on pre-defined styles, it can generalize to new styles with small training costs (a few minutes).

## 5.2 User Study

We conduct a user study to compare MM-NeRF with ARF [67] and StyleRF [34] on stylization quality and consistency. We chose three pairs of 3D scenes and styles to generate stylized novel views and videos with these stylization methods. For each pair of 3D scenes and styles, we ask the participant to vote for the best one in four indicators, including consistency across different views, overall quality, detail quality, and whether to be consistent with the color tone of styles. As shown in Table 3, MM-NeRF performs better than other methods, and over 70% of participants vote that MM-NeRF generates novel views with the most texture details.

## 5.3 Ablation Study

**Semantic Consistency of Multimodal Guidance.** To further verify that MM-NeRF can generate semantically consistent stylized images with multimodal styles, we organize comparative experiments on 360° scenes [65] and large-scale scenes [28]. As shown in Figure 6, MM-NeRF can ensure semantic consistency and generate similar stylized images with an image and text of the same style.

**Multi-Head Learning Scheme.** To verify the effectiveness of MLS, we compare MLS with a parameter predictor that all styles share a common prediction head. As shown in Figure 7 (b) and Figure 7 (d), MM-NeRF generates stylized results with clearer texture details. In addition, as shown in Figure 7 (c) and Figure 7 (e), for new styles, parameter predictor with single head fail to generalize to new styles and MLS can also generate high-quality results.

## 6 Conclusion

In this paper, we propose MM-NeRF, the first unified framework of multimodal-guided 3D multi-style transfer of NeRF. We also analyze the problem of lacking details in existing 3D stylization schemes and propose a novel Multi-Head Learning Scheme (MLS) to decompose the optimization difficulty caused by inconsistencies of multi-view and multi-style. Based on MLS, we propose a new incremental learning mechanism to generalize MM-NeRF to any new styles with small training costs. Experiments on three widely-used datasets demonstrate that MM-NeRF outperforms prior methods and exhibits advantages in 3D consistency, multimodal guidance, and semantic consistency of multimodal guidance.

## References

- [1] Metal albums artwork. <https://www.kaggle.com/datasets/benjamnmachn/metal-albums-artwork>.
- [2] Painters by numbers, wikiart. <https://www.kaggle.com/c/painter-by-numbers>.
- [3] Jean-Baptiste Alayrac, Adria Recasens, Rosalia Schneider, Relja Arandjelović, Jason Ramapuram, Jeffrey De Fauw, Lucas Smaira, Sander Dieleman, and Andrew Zisserman. Self-supervised multimodal versatile networks. *NeurIPS*, 33:25–37, 2020.
- [4] Jie An, Siyu Huang, Yibing Song, Dejing Dou, Wei Liu, and Jiebo Luo. Artflow: Unbiased image style transfer via reversible neural flows. In *CVPR*, pages 862–871, 2021.
- [5] 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*, pages 5855–5864, 2021.
- [6] Alexander Bergman, Petr Kellnhofer, Wang Yifan, Eric Chan, David Lindell, and Gordon Wetzstein. Generative neural articulated radiance fields. *NeurIPS*, 35:19900–19916, 2022.
- [7] Xu Cao, Weimin Wang, Katashi Nagao, and Ryosuke Nakamura. Psnet: A style transfer network for point cloud stylization on geometry and color. In *WACV*, pages 3337–3345, 2020.
- [8] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *CVPR*, pages 16123–16133, 2022.
- [9] Lele Chen, Sudhanshu Srivastava, Zhiyao Duan, and Chenliang Xu. Deep cross-modal audio-visual generation. In *ACM MM*, pages 349–357, 2017.
- [10] Xinghao Chen, Yiman Zhang, Yunhe Wang, Han Shu, Chunjing Xu, and Chang Xu. Optical flow distillation: Towards efficient and stable video style transfer. In *ECCV*, pages 614–630. Springer, 2020.
- [11] Yaosen Chen, Qi Yuan, Zhiqiang Li, Yuegen Liu, Wei 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.
- [12] Yuedong Chen, Qianyi Wu, Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. Sem2nerf: Converting single-view semantic masks to neural radiance fields. In *ECCV*, pages 730–748, 2022.
- [13] 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*, pages 1475–1484, 2022.

- [14] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197, 2020.
- [15] Yingying Deng, Fan Tang, Weiming Dong, Wen Sun, Feiyue Huang, and Changsheng Xu. Arbitrary style transfer via multi-adaptation network. In *ACM MM*, pages 2719–2727, 2020.
- [16] Yingying Deng, Fan Tang, Weiming Dong, Haibin Huang, Chongyang Ma, and Changsheng Xu. Arbitrary video style transfer via multi-channel correlation. In *AAAI*, volume 35, pages 1210–1217, 2021.
- [17] William Falcon and The PyTorch Lightning team. PyTorch Lightning, March 2019. URL <https://github.com/Lightning-AI/lightning>.
- [18] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *CVPR*, pages 8649–8658, 2021.
- [19] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR*, pages 2414–2423, 2016.
- [20] Simon Ging, Mohammadreza Zolfaghari, Hamed Pirsiavash, and Thomas Brox. Coot: Cooperative hierarchical transformer for video-text representation learning. *NeurIPS*, 33:22605–22618, 2020.
- [21] Yuan-Chen Guo, Di Kang, Linchao Bao, Yu He, and Song-Hai Zhang. Nerfren: Neural radiance fields with reflections. In *CVPR*, pages 18409–18418, 2022.
- [22] Wangli Hao, Zhaoxiang Zhang, and He Guan. Cmcgan: A uniform framework for cross-modal visual-audio mutual generation. In *AAAI*, volume 32, 2018.
- [23] Haozhi Huang, Hao Wang, Wenhan Luo, Lin Ma, Wenhao Jiang, Xiaolong Zhu, Zhifeng Li, and Wei Liu. Real-time neural style transfer for videos. In *CVPR*, pages 783–791, 2017.
- [24] Hsin-Ping Huang, Hung-Yu Tseng, Saurabh Saini, Maneesh Singh, and Ming-Hsuan Yang. Learning to stylize novel views. In *ICCV*, pages 13869–13878, 2021.
- [25] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, pages 1501–1510, 2017.
- [26] Yi-Hua Huang, Yue He, Yu-Jie Yuan, Yu-Kun Lai, and Lin Gao. Stylizednerf: consistent 3d scene stylization as stylized nerf via 2d-3d mutual learning. In *CVPR*, pages 18342–18352, 2022.
- [27] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. Neuman: Neural human radiance field from a single video. In *European Conference on Computer Vision*, pages 402–418. Springer, 2022.
- [28] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, pages 694–711. Springer, 2016.
- [29] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ToG*, 36(4):1–13, 2017.
- [30] Gang Li, Heliang Zheng, Chaoyue Wang, Chang Li, Changwen Zheng, and Dacheng Tao. 3ddesigner: Towards photorealistic 3d object generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2211.14108*, 2022.
- [31] Meng Li, Shangyin Gao, Yihui Feng, Yibo Shi, and Jing Wang. Content-oriented learned image compression. In *ECCV*, pages 632–647. Springer, 2022.
- [32] Xueting Li, Sifei Liu, Jan Kautz, and Ming-Hsuan Yang. Learning linear transformations for fast image and video style transfer. In *CVPR*, pages 3809–3817, 2019.
- [33] Yitong Li, Martin Min, Dinghan Shen, David Carlson, and Lawrence Carin. Video generation from text. In *AAAI*, volume 32, 2018.

- [34] Kunhao Liu, Fangneng Zhan, Yiwen Chen, Jiahui Zhang, Yingchen Yu, Abdulmoteleb El Saddik, Shijian Lu, and Eric Xing. Stylerf: Zero-shot 3d style transfer of neural radiance fields. In *CVPR*, 2023.
- [35] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, pages 7210–7219, 2021.
- [36] 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.
- [37] Fangzhou Mu, Jian Wang, Yicheng Wu, and Yin Li. 3d photo stylization: Learning to generate stylized novel views from a single image. In *CVPR*, pages 16273–16282, 2022.
- [38] Thomas Müller, Fabrice Rousselle, Jan Novák, and Alexander Keller. Real-time neural radiance caching for path tracing. *arXiv preprint arXiv:2106.12372*, 2021.
- [39] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ToG*, 41(4):1–15, 2022.
- [40] Thu Nguyen-Phuoc, Feng Liu, and Lei Xiao. Snerf: stylized neural implicit representations for 3d scenes. *ToG*, 41(4):1–11, 2022.
- [41] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *CVPR*, pages 11453–11464, 2021.
- [42] Yingwei Pan, Tao Mei, Ting Yao, Houqiang Li, and Yong Rui. Jointly modeling embedding and translation to bridge video and language. In *CVPR*, pages 4594–4602, 2016.
- [43] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, pages 5865–5874, 2021.
- [44] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- [45] Zhongwei Qiu, Qiansheng Yang, Jian Wang, and Dongmei Fu. Dynamic graph reasoning for multi-person 3d pose estimation. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 3521–3529, 2022.
- [46] Zhongwei Qiu, Qiansheng Yang, Jian Wang, and Dongmei Fu. Ivt: An end-to-end instance-guided video transformer for 3d pose estimation. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 6174–6182, 2022.
- [47] Zhongwei Qiu, Qiansheng Yang, Jian Wang, Haocheng Feng, Junyu Han, Errui Ding, Chang Xu, Dongmei Fu, and Jingdong Wang. Psvt: End-to-end multi-person 3d pose and shape estimation with progressive video transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21254–21263, 2023.
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021.
- [49] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022.
- [50] Manuel Ruder, Alexey Dosovitskiy, and Thomas Brox. Artistic style transfer for videos and spherical images. *IJCV*, 126(11):1199–1219, 2018.
- [51] Viktor Rudnev, Mohamed Elgharib, William Smith, Lingjie Liu, Vladislav Golyanik, and Christian Theobalt. Nerf for outdoor scene relighting. In *ECCV*, pages 615–631. Springer, 2022.



- [52] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 35:36479–36494, 2022.
- [53] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *CVPR*, pages 7495–7504, 2021.
- [54] Jan Svoboda, Asha Anoosheh, Christian Osendorfer, and Jonathan Masci. Two-stage peer-regularized feature recombination for arbitrary image style transfer. In *CVPR*, pages 13816–13825, 2020.
- [55] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, pages 402–419. Springer, 2020.
- [56] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>, 2022.
- [57] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *CVPR*, pages 3835–3844, 2022.
- [58] Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Nerf-art: Text-driven neural radiance fields stylization. *arXiv preprint arXiv:2212.08070*, 2022.
- [59] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021.
- [60] Wenjing Wang, Jizheng Xu, Li Zhang, Yue Wang, and Jiaying Liu. Consistent video style transfer via compound regularization. In *AAAI*, volume 34, pages 12233–12240, 2020.
- [61] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4): 600–612, 2004.
- [62] Tianyi Wei, Dongdong Chen, Wenbo Zhou, Jing Liao, Zhentao Tan, Lu Yuan, Weiming Zhang, and Nenghai Yu. Hairclip: Design your hair by text and reference image. In *CVPR*, pages 18072–18081, 2022.
- [63] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In *CVPR*, pages 16210–16220, 2022.
- [64] Jaejun Yoo, Youngjung Uh, Sanghyuk Chun, Byeongkyu Kang, and Jung-Woo Ha. Photorealistic style transfer via wavelet transforms. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9036–9045, 2019.
- [65] Kaan Yücer, Alexander Sorkine-Hornung, Oliver Wang, and Olga Sorkine-Hornung. Efficient 3d object segmentation from densely sampled light fields with applications to 3d reconstruction. *ToG*, 35(3):1–15, 2016.
- [66] Kai Zhang, Gernot Riegler, Noah Snaveley, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- [67] Kai Zhang, Nick Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snaveley. Arf: Artistic radiance fields. In *ECCV*, pages 717–733. Springer, 2022.
- [68] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pages 586–595, 2018.
- [69] Xiaoshuai Zhang, Sai Bi, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Nerfusion: Fusing radiance fields for large-scale scene reconstruction. In *CVPR*, pages 5449–5458, 2022.

- [70] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ToG*, 40(6):1–18, 2021.
- [71] Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara L Berg. Visual to sound: Generating natural sound for videos in the wild. In *CVPR*, pages 3550–3558, 2018.

## A Additional Implementation Details

### A.1 Network

To improve training and inference efficiency, Instant Neural Graphics Primitives (Instance-NGP) [39, 38, 64, 14, 61, 57] is employed to implement the basic NeRF and Ray Marching in MM-NeRF.

Spatial position  $\mathbf{x}$  and viewing direction  $\mathbf{d}$  are encoded based on the trainable multi-resolution grid [39] and the multi-resolution sequence of spherical harmonics, respectively. Our implementation is based on PyTorch-Lightning [17]. Stable Diffusion [49] is implemented with Diffusers of HuggingFace [56].

### A.2 Pre-Training of Multi-Head Learning Scheme

We have observed that randomly initializing Multi-Head Learning Scheme (MLS) leads to unstable stylization training. Before stylization training, we introduce a pre-training process of MLS to fix this issue. At each iteration, we randomly sample a style from the style set  $S$  and predict the parameters of the color head of NeRF. MLS is optimized to minimize Mean Square Error (MSE) between predicted parameters and pre-trained parameters of the color head of NeRF:

$$\mathcal{L}_p = \sum_{i=1}^{N_s} \|\hat{\mathbf{w}}_{s^{(i)}} - \mathbf{w}\|_2^2, \quad (10)$$

where  $N_s$  is the total number of styles in  $S$ ,  $\hat{\mathbf{w}}_{s^{(i)}}$  is the predicted parameters with  $s^{(i)}$  by MLS and  $\mathbf{w}$  is the pre-trained parameters of the color head of NeRF.

**Analysis.** Stylization results are sensitive to NeRF parameters, which means that MLS has a complex parameter space. Minor changes in MLS predictions may also lead to drastic changes in stylized results. Therefore, randomly initialized MLS causes extreme instability. The above issue can be solved by reducing the learning rate, but reducing the learning rate leads to a decrease in convergence speed. Based on the above analysis, introducing the pre-training process of MLS to ensure a good initialization is the most appropriate method to stabilize the stylization training process.

## B Calculation of Metrics

We render a video of the stylized scene with a moving camera for each pair of a scene and a style. The optical flow of this video is computed by RAFT [55] and Temporal Warping Error (TWE) is computed as follows:

$$E_{TWE}(I_g^{(i)}, I_g^{(j)}) = MSE(I_g^{(i)}, M^{(i,j)}, W^{(i,j)}(I_g^{(j)})), \quad (11)$$

where  $I_g^{(i)}$  and  $I_g^{(j)}$  are two stylized views in this video,  $M^{(i,j)}$  is the occlusion mask and  $W^{(i,j)}$  is the warping function.  $I_g^{(j)}$  is warped to view  $i$  and TWE is the masked MSE between  $I_g^{(i)}$  and  $W^{(i,j)}(I_g^{(j)})$ . Similarly, the warped LPIPS is the masked LPIPS between  $I_g^{(i)}$  and  $W^{(i,j)}(I_g^{(j)})$ :

$$E_{TWE}(I_g^{(i)}, I_g^{(j)}) = LPIPS(I_g^{(i)}, M^{(i,j)}, W^{(i,j)}(I_g^{(j)})), \quad (12)$$

Huang et al. [24] use  $M^{(i,j)}$  to choose valid pixels. These pixels are used to compute the ‘‘spatial average’’ in [68]. As LPIPS is a type of feature loss, it is not accurate enough to compute the masked LPIPS with masks of pixels and we compute it as follows:

$$E_{LPIPS}(I_g^{(i)}, I_g^{(j)}) = LPIPS(I_g^{(i)}, W^{(i,j)}(I_g^{(j)})'), \quad (13)$$

where  $W^{(i,j)}(I_g^{(j)})'$  is the replaced reconstructed image [31] of  $W^{(i,j)}(I_g^{(j)})$ :

$$W^{(i,j)}(I_g^{(j)})' = (1 - M^{(i,j)}) \circ I_g^{(i)} + M^{(i,j)} \circ W^{(i,j)}(I_g^{(j)}). \quad (14)$$

Specifically, short-range consistency is measured by  $I_g^{(i)}$  and  $I_g^{(i+1)}$ , which are the two adjacent images in the video, and long-range consistency is measured by  $I_g^{(i)}$  and  $I_g^{(i+7)}$ , which are two images separated by seven frames in the video.

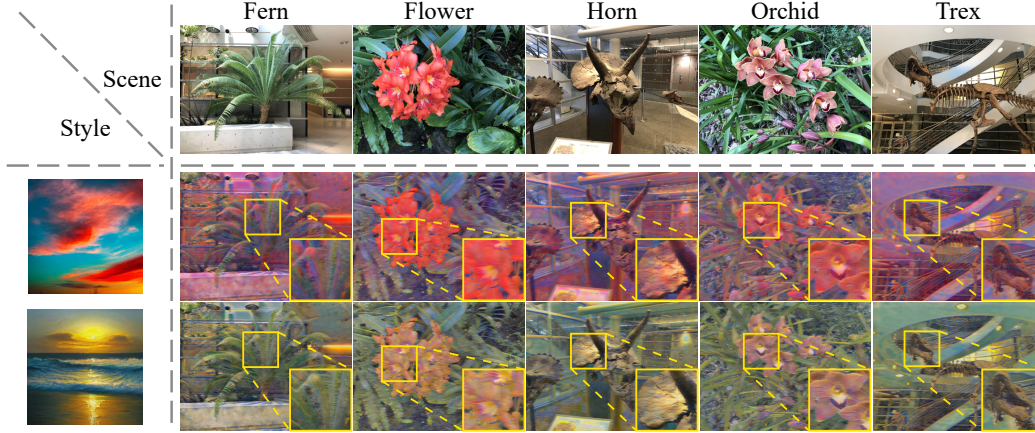


Figure 8: Failure Cases. MM-NeRF tends to generate blurry results with styles lacking brushstrokes.

## C Limitations

Although NeRF can model forward-facing scenes [36] well, its modeling accuracy still needs to be improved for large-scale scenes [29]. As MM-NeRF is based on NeRF, when NeRF fails to model scenes with texture details, MM-NeRF can not achieve high-quality stylized novel views. In addition, the 2D style transfer model, which is used to generate stylized images to optimize MM-NeRF, also impacts the performance of MM-NeRF. In MM-NeRF, AdaIN [25] is employed. When the strokes of the reference style are not obvious, AdaIN can not generate stylized results with texture details. As shown in Figure 8, MM-NeRF tends to generate blurry results in this case.