

NeRF-Art: Text-Driven Neural Radiance Fields Stylization

CAN WANG, City University of Hong Kong
 RUIXIANG JIANG, The Hong Kong Polytechnic University
 MENGLEI CHAI, Snap Inc.
 MINGMING HE, Netflix
 DONGDONG CHEN, Microsoft Cloud AI
 JING LIAO*, City University of Hong Kong



Fig. 1. **NeRF-Art Results.** Our *NeRF-Art* stylizes a pre-trained NeRF to match the desired style described by a text prompt. It modulates not only appearance but also geometry of NeRF.

As a powerful representation of 3D scenes, the neural radiance field (NeRF) enables high-quality novel view synthesis from multi-view images. Stylizing NeRF, however, remains challenging, especially on simulating a text-guided style with both the appearance and the geometry altered simultaneously. In this paper, we present *NeRF-Art*, a text-guided NeRF stylization approach that manipulates the style of a pre-trained NeRF model with a simple text prompt. Unlike previous approaches that either lack sufficient geometry deformations and texture details or require meshes to guide the stylization, our method can shift a 3D scene to the target style characterized by desired geometry and appearance variations without any mesh guidance. This is achieved by introducing a novel global-local contrastive learning strategy, combined with the directional constraint to simultaneously control both the trajectory and the strength of the target style. Moreover, we adopt a weight regularization

method to effectively suppress cloudy artifacts and geometry noises which arise easily when the density field is transformed during geometry stylization. Through extensive experiments on various styles, we demonstrate that our method is effective and robust regarding both single-view stylization quality and cross-view consistency. The code and more results can be found in our project page: <https://cassiepython.github.io/nerfart/>.

CCS Concepts: • **Computing methodologies** → **Computer graphics**.

Additional Key Words and Phrases: Stylization, Neural Radiance Fields, CLIP

1 INTRODUCTION

Artistic works depict the world in various creative and imaginative styles, evolving along with human progress. While primarily driven by professionals, the generation of artistic content is now more accessible to average users than ever before, empowered by the recent research on visual artistic stylization. In the era of deep learning, technical advances are gradually reshaping how people

*Corresponding Author.

Authors' addresses: Can Wang, City University of Hong Kong, cwang355-c@my.cityu.edu.hk; Ruixiang Jiang, The Hong Kong Polytechnic University, cwang355-c@my.cityu.edu.hk; Menglei Chai, Snap Inc. cmlatsim@gmail.com; Mingming He, Netflix, hmm.lillian@gmail.com; Dongdong Chen, Microsoft Cloud AI, cddyf@gmail.com; Jing Liao*, City University of Hong Kong, jingliao@cityu.edu.hk.

create, consume, and share art, from real-time entertainment to concept design. Since neural style transfer [Chen et al. 2017b; Gatys et al. 2016; Sheng et al. 2018; Shu et al. 2021; Zhao et al. 2014] shows the potential of encoding and changing visual styles via deep neural networks, a significant amount of effort has been devoted to effectively and efficiently migrating the style of an arbitrary image [Gatys et al. 2016; Huang and Belongie 2017; Li et al. 2017; Liao et al. 2017] or a specific domain [Lee et al. 2020; Zhu et al. 2017] to the content image. Despite impressive results, these methods are limited to stylizing a single view captured by the content image.

Motivated by the increasing demand for 3D asset creation, our goal is to stylize *3D content* from *multi-view input*, in contrast to single-image stylization. In the domain of 3D representation, previous methods typically take explicit models (e.g., meshes [Han et al. 2021a; Höllein et al. 2021; Kato et al. 2018; Ye et al. 2021; Zhang et al. 2020a], voxels [Guo et al. 2021; Klehm et al. 2014], and point clouds [Cao et al. 2020; Lin et al. 2018]) followed by differentiable rendering for multi-view stylization. These methods enable intuitive control over the geometry but suffer from the limited capacity for modeling and rendering complex scenes. The recent implicit representation of neural radiance field (NeRF) [Deng et al. 2022; Mildenhall et al. 2020; Wang et al. 2022; Yang et al. 2022; Zhang et al. 2022b] significantly improves the quality of novel view synthesis, satisfying our needs for a general representation of various scenes and objects. However, while enjoying the superior scene reconstruction quality of NeRF, the curse of its highly implicit volumetric representation of appearance and geometry, parameterized and entangled by dense MLP networks, makes NeRF more challenging to stylize through jointly transforming the encoded color and shape.

Very recently, pioneering NeRF stylization works [Chiang et al. 2022; Fan et al. 2022] have made exhilarating progress on appearance style transfer of 3D scenes. However, their style guidance is limited to image reference, which, although being adopted as one common way to specify the target style, is not always a perfect solution for every scenario—obtaining appropriate style images that both reflect the target style and match the source content might not be easy or even possible in many cases. Therefore, finding another simple, natural, and expressive form of guidance becomes an attractive idea. Thanks to the parallel advances in language-vision models, stylization with natural language is no longer a fantasy. As demonstrated by recent text-guided stylization works [Gal et al. 2021; Hong et al. 2022; Michel et al. 2021; Wei et al. 2021], compared to image-guided approaches, short text prompts provide 1) an extremely intuitive and user-friendly way to specify styles, 2) a flexible control over various styles from abstract ones like a certain concept to very concrete ones like a famous painting or character, and 3) a view-independent representation that is free from content alignment and naturally benefits cross-view consistency.

Yet, with the existing approaches, it is still challenging to stylize the implicit representation of NeRF via a simple text prompt. Learning a latent space helps constrain the geometry and texture modulations [Wang et al. 2021a], but it is often data-dependent and laborious. Some efforts directly enforce style directions (Figure 3) between the rendered views of NeRF and the text in the CLIP [Radford et al. 2021] embedding space. In addition, background augmentation [Jain et al. 2022] and mesh guidance [Hong et al. 2022] have

been proposed to improve the geometry and texture modulations. However, they still suffer from insufficient geometry deformations and texture details.

In this work, we propose *NeRF-Art*, a new text-driven NeRF stylization method. Given a pre-trained NeRF model and a single text prompt, our method enables consistent novel view synthesis with both appearance and geometry transformed, adhering to the specified style. This is achieved by combining the recent large-scale Language-Vision model (i.e., CLIP) with NeRF, which is non-trivial due to several challenges. Directly applying the supervision from CLIP to NeRF by constraining the similarity between the rendered views and the text in the embedding space as [Gal et al. 2021] is insufficient to ensure the desired style strength. To tackle this problem, we design a CLIP-based contrastive loss to properly strengthen the stylization, by bringing the results closer to the target style and farther away from other styles pre-defined as negative samples. To further ensure the uniformity of the style over the whole scene, we extend our contrastive constraint to a hybrid global-local framework to cover both global structures and local details. In addition, to support geometry stylization jointly with appearance, we relax the constraints on the density of the pre-trained NeRF and adopt a weight regularization to effectively reduce cloudy artifacts and geometry noises when altering the density field. In experiments, we first evaluate text description selection for stylization and then test our method on various styles and demonstrate text guidance’s effectiveness and flexibility for NeRF stylization. Furthermore, we conduct a user study to show that our method achieves the best visual-pleasing results compared to related methods. We also extract the mesh from the stylized NeRF to show the geometry modulation ability of our method and integrate with different baselines to demonstrate the generalization ability of our method to various NeRF-like models.

2 RELATED WORK

Neural Style Transfer on Images and Videos. Artistic image stylization is a long-standing research area. Traditional methods use handcrafted features to simulate styles [Hertzmann 1998; Hertzmann et al. 2001]. With the fast development of deep learning, neural networks have been applied to style transfer from either an arbitrary image [Gatys et al. 2016; Huang and Belongie 2017; Johnson et al. 2016; Kolkin et al. 2019; Li et al. 2017; Liao et al. 2017] or a specific domain [Huang et al. 2021a, 2018; Lee et al. 2020; Zhu et al. 2017], and achieved impressive results. By enforcing temporal smoothness constraints defined on optical flows, neural style transfer has been successfully extended to videos [Chen et al. 2017a, 2020; Ruder et al. 2016]. However, both image and video stylization methods are restricted to the given views. Simply combining the neural style transfer and novel view synthesis methods without considering 3D geometry will lead to blurriness or view inconsistencies.

Neural Stylization on Explicit 3D Representations. With the increasing demand for 3D content, neural style transfer has been extended to explicit 3D representations. The work [Chen et al. 2018] first considers the cross-view disparity consistency and applies style transfer on stereoscopic images or videos. Later, considering the voxel is the most compatible representation for CNNs, *SKPN* [Guo

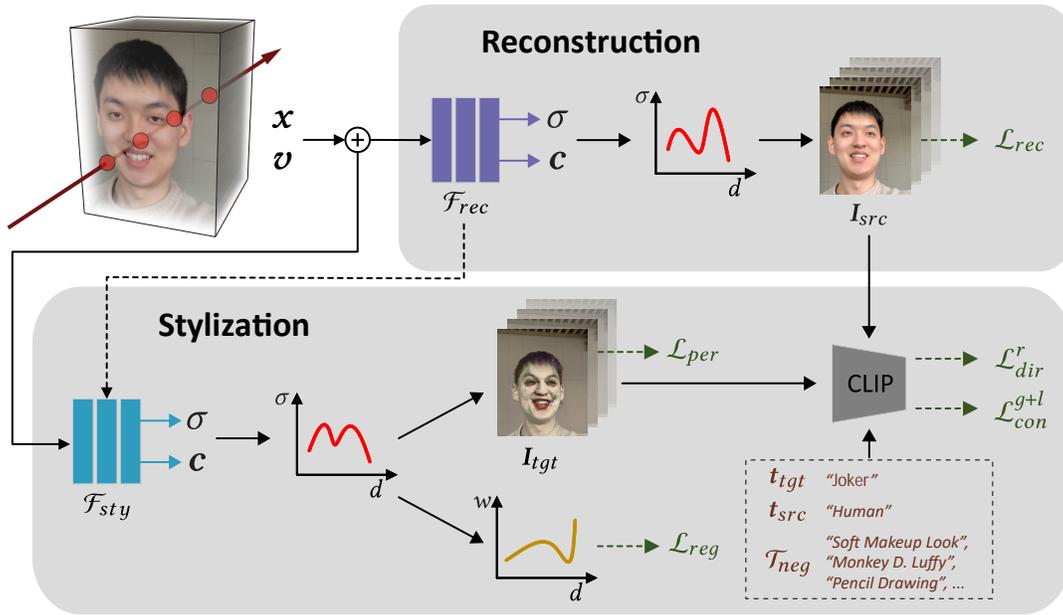


Fig. 2. **NeRF-Art Pipeline.** In the *reconstruction* stage, our method first pre-trains the NeRF model \mathcal{F}_{rec} of the target scene from multi-view input with reconstruction loss \mathcal{L}_{rec} . In the *stylization* stage, our method stylized NeRF model \mathcal{F}_{rec} to \mathcal{F}_{sty} , guided by a text prompt t_{tgt} , using a combination of relative directional loss \mathcal{L}_{dir}^r and global-local contrastive loss \mathcal{L}_{con}^{g+l} in the CLIP embedding space, plus weight regularization loss \mathcal{L}_{reg} and perceptual loss \mathcal{L}_{per} .

et al. 2021] encodes volume using convolutional blocks and stylizes it by deep features extracted from a reference image. As for mesh stylization, differential rendering allows for backpropagating style transfer objectives from rendered images to 3D meshes. According to whether the geometry or texture are allowed to be optimized, existing mesh style transfer methods achieve three different effects: texture stylization [Höller et al. 2021; Mordvintsev et al. 2018], geometric stylization [Liu et al. 2018], and joint stylization [Han et al. 2021b; Kato et al. 2018; Yin et al. 2021]. Another line of work uses point clouds as the 3D proxy to guarantee 3D consistency in stylizing novel views from either a single image [Mu et al. 2021] or multiple frames [Huang et al. 2021b]. In these works, point-wise features extracted from pre-trained PointNet [Qi et al. 2017] or GCN [Li et al. 2021a] are stylized by feature transform algorithms, e.g., adaptive normalization, and then rendered to novel views. Despite the successes, these 3D stylization methods are difficult to generalize to complicated objects or scenes with dedicated structures, limited by the expressiveness of explicit 3D representations.

Neural Stylization on NeRF. To address the inherent limitations of explicit representations, implicit methods have recently received much attention. NeRF is a seminal one that is able to represent complex scenes by parameterizing the implicit function as MLP networks. A large number of follow-up works are presented to improve its efficiency [Deng et al. 2021; Garbin et al. 2021; Lindell et al. 2021; Müller et al. 2022; Reiser et al. 2021; Yu et al. 2021a], quality [Arandjelović and Zisserman 2021; Barron et al. 2021; Ma et al. 2021; Zhang et al. 2020b], controllability [Liu et al. 2021; Srinivasan et al. 2021; Wang et al. 2021a; Zhang et al. 2021], and generalization [Gao et al. 2021; Jain et al. 2021; Li et al. 2021b; Niemeyer et al. 2021; Noguchi

et al. 2021; Park et al. 2021a,b; Peng et al. 2021; Pumarola et al. 2021; Tretschk et al. 2021; Xian et al. 2021; Yu et al. 2021b]. Inspired by the power of NeRF, three very recent works [Chiang et al. 2022; Huang et al. 2022; Zhang et al. 2022a] adopt it for 3D stylization. They design the stylization network to predict color-related parameters in the NeRF model based on a reference style. And the stylization network is trained either by imposing the image style transfer losses [Gatys et al. 2016; Zhang et al. 2022a] on rendered views [Chiang et al. 2022] or being supervised by a mutually learnt image stylization network [Huang et al. 2022]. These works have achieved consistent results in novel-view stylization. However, their stylization is still restricted to appearance only because they do not adjust density parameters in the NeRF model. In contrast, our method supports both appearance and geometric stylization to better mimic the reference style. Moreover, they rely on reference images for stylization, while we seek to stylize the scenes via simple text prompts.

Text-Driven Stylization. Compared to image references, a natural language prompt is a more intuitive and user-friendly way to specify the style. Therefore, a current line of works shifted away from image reference towards text guidance, with the help of the pre-trained CLIP [Radford et al. 2021], which bridges texts and images by jointly learning a shared latent space. The pioneering work *StyleGAN-NADA* [Gal et al. 2021] proposes a directional CLIP loss for transferring the pre-trained StyleGAN2 model [Karras et al. 2020] to the target domain with the desired style described by a textual prompt. However, it is an image-based method and will lead to inconsistencies when applied to stylizing multiple views. In the 3D world, *Text2Mesh* [Michel et al. 2021] uses CLIP to guide the stylization of a given 3D mesh by learning a displacement map for

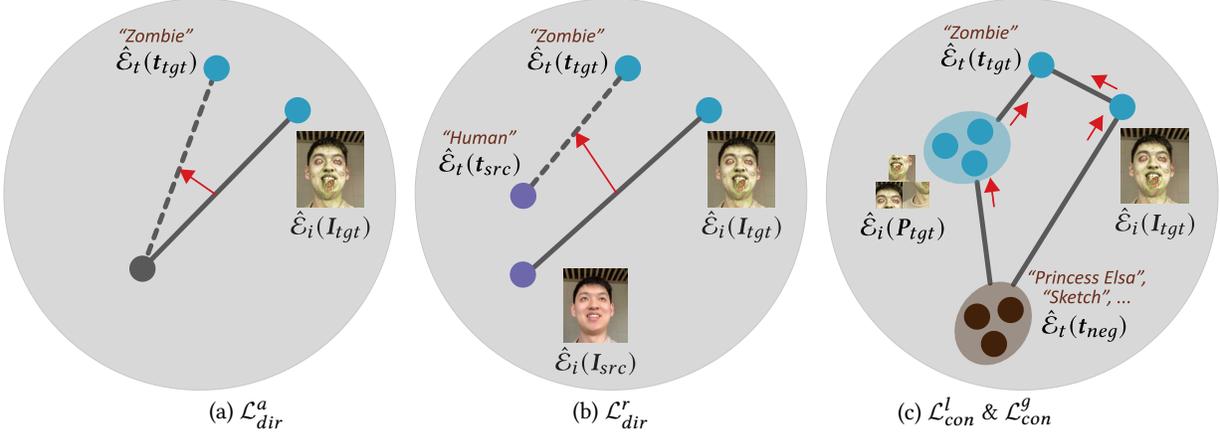


Fig. 3. **CLIP-Guided Stylization Losses.** (a) The absolute directional loss; (b) The relative directional loss; (c) The global and local contrastive loss.

geometry deformation and vertex colors for texture stylization. The contemporary work *AvatarCLIP* [Hong et al. 2022] further supports driving a stylized human mesh using natural languages. Despite their success, these methods are limited to mesh input. In contrast, our method is able to stylize 3D scenes with better visual quality and view consistency without any mesh input.

3 OVERVIEW

As illustrated in Figure 2, our approach is simply decomposed into reconstruction and stylization stages. In what follows, after briefly reviewing our 3D photography representation with NeRF (§ 3.1), we focus on introducing our text-guided stylization method. Specifically, we first formulate the directional CLIP loss for stylization, which leverages the power of the pre-trained Language-Vision model (§ 4.1). Then, we introduce our global-local contrastive learning framework to cope with the stylization strength issue of the directional CLIP loss (§ 4.2). Next, we introduce a weight regularization term to alleviate the cloudy artifacts caused and geometry noises by the stylization process (§ 4.3). Finally, we conclude this section with the overall training strategy of the entire pipeline (§ 4.4).

3.1 Preliminary on NeRF Scene Representation

We take NeRF as our 3D scene representation, which defines a continuous volumetric field as implicit functions, parameterized by MLP networks \mathcal{F} . Given a single spatial coordinate $\mathbf{x} = (x, y, z)$ and its corresponding view direction $\mathbf{d} = (\phi, \theta)$, the network predicts the density σ and view-dependent radiance $\mathbf{c} = (r, g, b)$, leading to the final color $C(\mathbf{r})$ of the camera ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ by accumulating K sample points along it, given the target view:

$$C(\mathbf{r}) = \sum_{k=1}^K T_k (1 - \omega_k) \mathbf{c}_k, \quad (1)$$

where $\omega_k = \exp(-\sigma_k(d_{k+1} - d_k))$ represents the transmittance of the ray segment $(k, k+1)$ and $T_k = \prod_{i=1}^{k-1} \omega_i$ is the accumulated transmittance from the origin to the sample k .

To train NeRF from a set of multi-view photos, a simple supervised reconstruction loss is adopted between the ground-truth pixel colors

$\hat{C}(\mathbf{r})$ from the training view and the NeRF prediction $C(\mathbf{r})$:

$$\mathcal{L}_{rec} = \sum_{\mathbf{r}} \|C(\mathbf{r}) - \hat{C}(\mathbf{r})\|_2^2. \quad (2)$$

4 TEXT-GUIDED NERF STYLIZATION

After optimizing the reconstructed NeRF model \mathcal{F}_{rec} from the multi-view input (§ 3.1), our goal is to train a stylized NeRF model \mathcal{F}_{sty} , which satisfies the style control of the target text prompt \mathbf{t}_{tgt} while preserving the content from \mathcal{F}_{rec} (Figure 2).

The CLIP model aligns the semantics of image and text in a joint embedding space, by utilizing the image encoder $\hat{\mathcal{E}}_i(\cdot)$ and the text encoder $\hat{\mathcal{E}}_t(\cdot)$. The semantic power of CLIP bridges the gap between natural language prompts and synthesized image pixels, making it possible to stylize NeRF scenes with text controls.

However, even with the powerful embedding space of CLIP, it remains challenging to achieve text-guided NeRF stylization that 1) preserves the original content from being washed away by the new style, 2) reaches the target style with proper strength that satisfies the semantics of the input text prompt, and 3) maintains cross-view consistency and avoids artifacts in the final NeRF model.

4.1 Trajectory Control w/ Directional CLIP Loss

An intuitive strategy for text-guided NeRF stylization would be to enforce the trajectory of the stylization in the CLIP space with an *absolute* directional CLIP loss that measures the cosine similarity $\langle \cdot, \cdot \rangle$ between the stylized NeRF rendering I_{tgt} and the target text prompt \mathbf{t}_{tgt} (Figure 3(a)):

$$\mathcal{L}_{dir}^a = \sum_{I_{tgt}} \left[1 - \langle \hat{\mathcal{E}}_i(I_{tgt}), \hat{\mathcal{E}}_t(\mathbf{t}_{tgt}) \rangle \right], \quad (3)$$

which guides NeRF rendering with a global direction of the target text, not depending on any reference starting point. This loss is first designed in *StyleCLIP* [Patashnik et al. 2021] to guide face image editing and further extended to generative NeRF editing in *CLIP-NeRF* [Wang et al. 2021a].

However, as observed in *StyleGAN-NADA* [Gal et al. 2021], this global loss could easily mode-collapse the generator and hurt the generation diversity of stylization. Therefore, a *relative* directional

loss is proposed, which transfers the source image I_{src} to the target domain guided by the CLIP-space trajectory embedded by a pair of text prompts ($\mathbf{t}_{src}, \mathbf{t}_{tgt}$) instead of a single one (Figure 3(b)). This relative directional CLIP loss for our NeRF stylization is defined as:

$$\mathcal{L}_{dir}^r = \sum_{I_{tgt}} \left[1 - \langle \hat{\mathcal{E}}_i(I_{tgt}) - \hat{\mathcal{E}}_i(I_{src}), \hat{\mathcal{E}}_t(\mathbf{t}_{tgt}) - \hat{\mathcal{E}}_t(\mathbf{t}_{src}) \rangle \right]. \quad (4)$$

Different from the single-image setting of *StyleGAN-NADA*, here, the training target I_{tgt} stands for an arbitrarily sampled view rendered by the stylized NeRF of the same scene, and the source image I_{src} is produced by the pre-trained NeRF model and shares the identical view as I_{tgt} . We will follow this convention hereinafter.

4.2 Strength Control w/ Global Contrastive Learning

As the directional CLIP loss (Equation (4)) works by measuring the similarity between the normalized unit directions of the embedded vectors, it can enforce the relative stylization trajectory. However, it struggles with preserving enough stylization strength in altering the pre-trained NeRF model.

To address this issue, we propose a contrastive learning strategy to control the stylization strength (Figure 3(c)). Specifically, in the framework of contrastive learning, with the rendered view I_{tgt} as the query target, we set positive samples to the target text prompt \mathbf{t}_{tgt} with the desired style and construct negative samples $\mathbf{t}_{neg} \in \mathcal{T}_{neg}$ by sampling a set of text prompts semantically irrelevant to I_{tgt} . In general, our contrastive loss in the CLIP space is defined as:

$$\mathcal{L}_{con} = - \sum_{I_{tgt}} \log \left[\frac{\exp(\mathbf{v} \cdot \mathbf{v}^+ / \tau)}{\exp(\mathbf{v} \cdot \mathbf{v}^+ / \tau) + \sum_{\mathbf{v}^-} \exp(\mathbf{v} \cdot \mathbf{v}^- / \tau)} \right], \quad (5)$$

where $\{\mathbf{v}, \mathbf{v}^+, \mathbf{v}^-\}$ are query, positive sample, and negative sample, respectively, and temperature τ is set to 0.07 in all our experiments. When defining the loss globally by treating the entire view I_{tgt} as the query anchor, we have the global contrastive loss \mathcal{L}_{con}^g with $\{\mathbf{v} = \hat{\mathcal{E}}_i(I_{tgt}), \mathbf{v}^+ = \hat{\mathcal{E}}_t(\mathbf{t}_{tgt}), \mathbf{v}^- = \hat{\mathcal{E}}_t(\mathbf{t}_{neg})\}$.

Ideally, this global contrastive loss cooperates with the directional CLIP loss, where the former defines the style trajectory that aligns with the target text, and the latter, at the same time, ensures the proper stylization magnitude by pushing along the style trajectory. However, the global contrastive loss still has trouble achieving sufficient and uniform stylization on the entire NeRF scene, leading to excessive stylization on certain parts and insufficient stylization in other regions. This may be attributed to the fact that CLIP focuses more attention on local regions with distinguishable features than the entire scene. Thus, this global contrastive loss can deliver a small value even when the overall stylization is insufficient or non-uniform. To achieve a more sufficient and balanced stylization, enforced by a more locally-attended contrastive learning approach, inspired by *PatchNCE* loss [Park et al. 2020], we propose a complementary local contrastive loss \mathcal{L}_{con}^l which sets queries to random local patches P_{tgt} cropped from I_{tgt} : $\{\mathbf{v} = \hat{\mathcal{E}}_i(P_{tgt}), \mathbf{v}^+ = \hat{\mathcal{E}}_t(\mathbf{t}_{tgt}), \mathbf{v}^- = \hat{\mathcal{E}}_t(\mathbf{t}_{neg})\}$.

Overall, we combine the global and local terms as our final global-local contrastive loss:

$$\mathcal{L}_{con}^{g+l} = \lambda_g \mathcal{L}_{con}^g + \lambda_l \mathcal{L}_{con}^l. \quad (6)$$

4.3 Artifact Suppression w/ Weight Regularization

Our pipeline aims to change not only the color but also the density of the pre-trained NeRF to achieve a joint stylization of appearance and geometry. However, allowing the training process to alter the density may lead to cloud-like semi-transparent artifacts near the camera and geometry noises, even if the pre-trained NeRF is perfectly clean. To alleviate that, we adopt a weight regularization loss to suppress geometric noises and encourage a more concentrated density distribution that better resembles real-world scenes.

Based on our NeRF notations (Equation (1)), weight of each ray sample is defined as the contribution to the final ray color: $w_k = T_k(1 - \omega_k)$, where $\sum_k w_k \leq 1$. Similar to the distortion loss in *mip-NeRF 360* [Barron et al. 2022], the weight regularization loss is defined as:

$$\mathcal{L}_{reg} = \sum_{I_{tgt}} \sum_{\mathbf{r}} \sum_{(i,j) \in K} w_i w_j \|d_i - d_j\|, \quad (7)$$

where for each ray \mathbf{r} of a randomly sampled view I_{tgt} , pairs of samples (i, j) with distances $\|d_i - d_j\|$ are sampled. But different from *mip-NeRF 360* [Barron et al. 2022] that optimizes the distances, we penalize those pairs with scattered large weights to suppress noise peaks and aggregate weights to the correct object surface.

4.4 Training

During training, we finetune the pre-trained NeRF model for stylization. The overall objective consists of three parts: text-guided stylization losses (including directional CLIP loss and global-local contrastive loss to control style trajectory and strength, respectively), content-preservation loss (we adopt VGG-based perceptual loss), and artifact suppression regularization loss:

$$\mathcal{L} = (\mathcal{L}_{dir}^r + \mathcal{L}_{con}^{g+l}) + \lambda_p \mathcal{L}_{per} + \lambda_r \mathcal{L}_{reg}. \quad (8)$$

Here we define the perceptual loss \mathcal{L}_{per} between the original and stylized NeRF renderings on certain pre-defined VGG layers $\psi \in \Psi$:

$$\mathcal{L}_{per} = \sum_{I_{tgt}} \sum_{\psi \in \Psi} \|\psi(I_{tgt}) - \psi(I_{src})\|_2^2. \quad (9)$$

It’s practically infeasible to train stylization on all rays due to backward gradient propagation’s prohibitively huge memory consumption. To address this issue, previous works either sample sparse rays to obtain coarse images or patches [Chiang et al. 2022; Hong et al. 2022; Jain et al. 2021; Schwarz et al. 2020] or render all rays to low resolution and then upsample with CNN networks [Niemeyer and Geiger 2021]. However, coarse renderings or patches lose style details and semantic structures, while upsampling harms the cross-view consistency. Instead, we adopt a much easier solution, which first renders all rays to obtain the whole image of an arbitrary view, calculates the stylization loss gradients in the forward process, and then back-propagates the gradients through NeRF at the patch level. This significantly reduces memory consumption and allows rendering high-resolution images for better stylization training.

5 EXPERIMENTS

5.1 Implementation Details

We implement our framework using *Pytorch*. In the reconstruction training stage, we sample 192 points for each ray and train our model for 6 epochs. We set the learning rate as 0.0005 and adopt

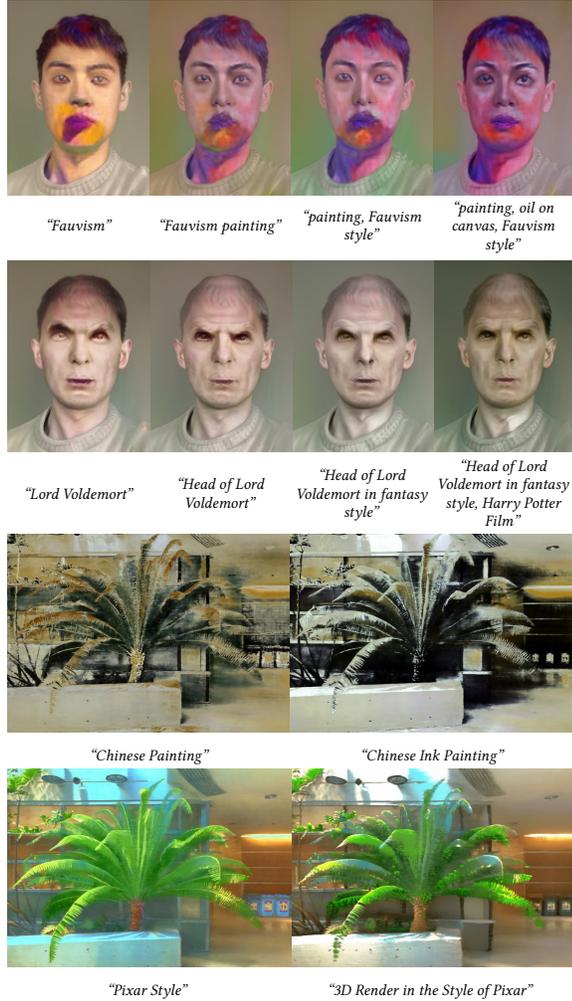


Fig. 4. **Text Evaluation.** We present descriptions at different detail levels for a specific style.

the Adam optimizer. While in the stylization training stage, we train our model for 4 epochs with the learning rate of 0.001 and use the *Adam* optimizer. We set hyper-parameters λ_g , λ_l , λ_p and λ_r as 0.2, 0.1, 2.0, and 0.1, respectively. To construct the negative samples, we manually collect around 200 text descriptions from *Pinterest* website, describing various styles, like “*Zombie*”, “*Tolkien elf*”, and “*Self-Portrait by Van Gogh*”. We set the patch size as the 1/10 of the original input in the local contrastive loss. Without loss of generality, we adopt VolSDF [Yariv et al. 2021] as the basic NeRF model for stylization.

5.2 Data Collection

Three self-portrait datasets are gathered under an in-the-wild condition by asking three users to capture selfies video for around 10 seconds with the front-facing camera. We finally received six video clips in around 10 seconds. After collecting these video clips under different views and expressions, we extract 100 frames for each

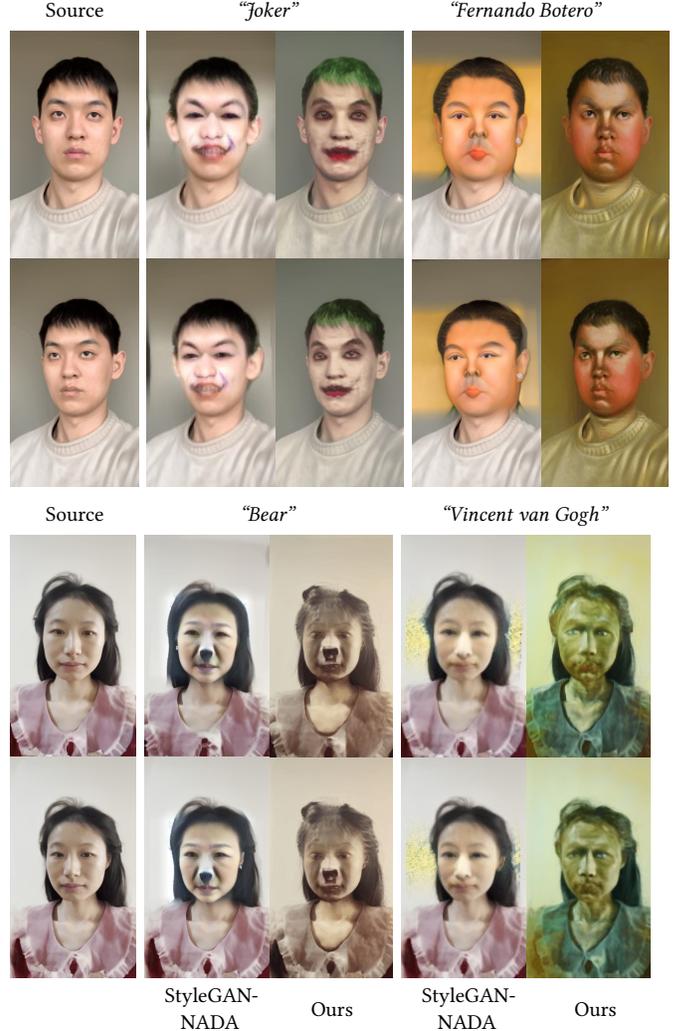


Fig. 5. **Comparisons.** Comparisons with the text-guided image stylization method *StyleGAN-NADA* [Gal et al. 2021].

video clip using Ffmpeg with 15 fps. Then these frames are resized to 270×480 . Then we estimate camera poses for these frames using COLMAP [Schonberger and Frahm 2016] with rigid relative camera pose constraints. We suppose frames in a video share the same intrinsics. We also reconstruct a lady from the H3DS dataset [Ramon et al. 2021]. We remove noise frames and obtain 31 sparse views. Moreover, we use the image size with 256×256 for stylization. We also adopt the Local Light Field Fusion (LLFF) dataset [Mildenhall et al. 2019] to stylize non-face scenes. LLFF dataset is composed of forward-facing scenes, with around 20 to 60 images.

5.3 Text Evaluation

As CLIP [Park et al. 2020] is sensitive to text prompts, we conduct a text description evaluation in Figure 4. When a text description refers to a style in general, not anyone in particular, the stylization

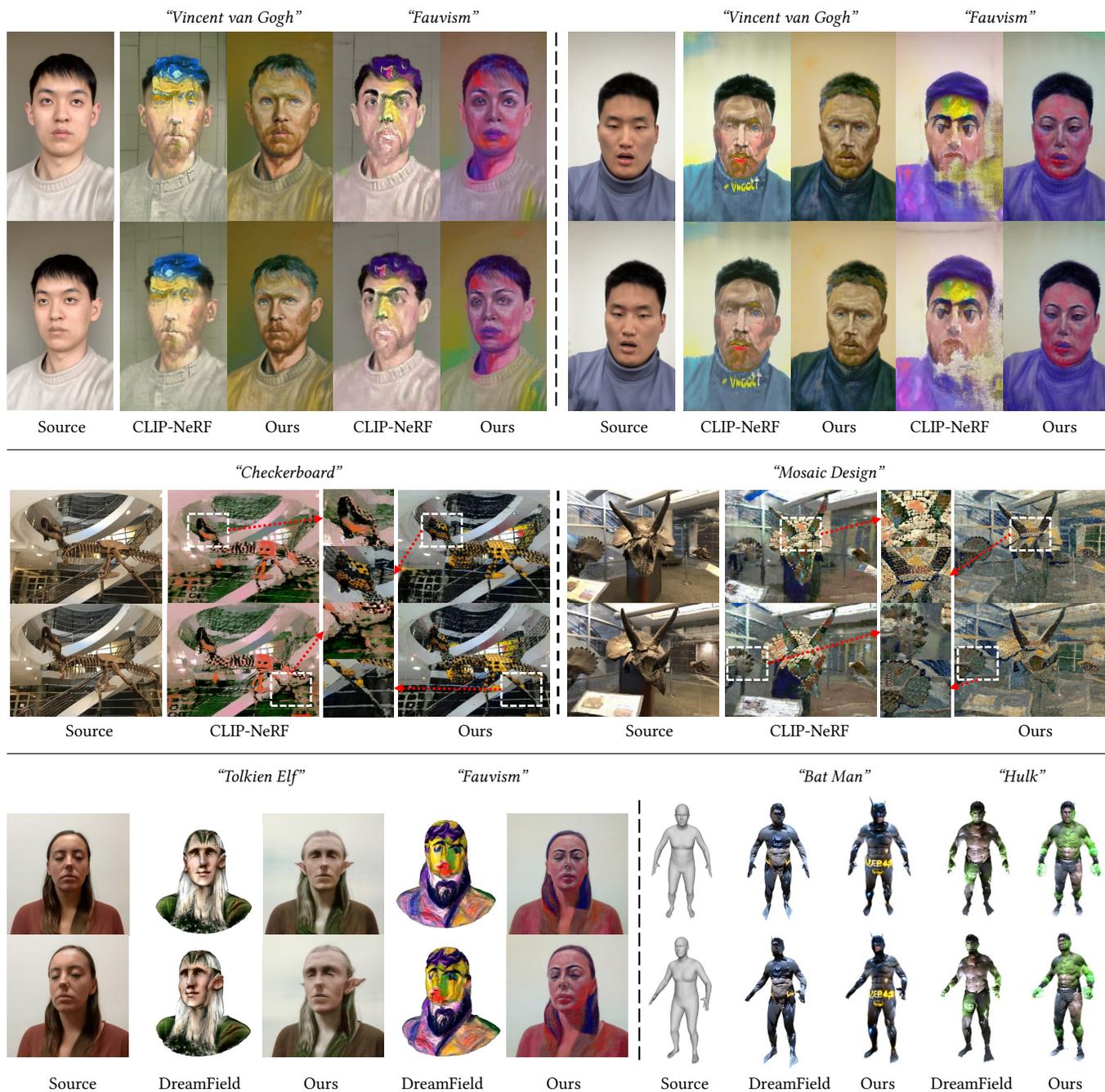


Fig. 6. **Comparisons.** Comparisons to text-guided NeRF stylization method *CLIP-NeRF* [Wang et al. 2021a] and *DreamField* [Jain et al. 2022].

can be insufficient. For example, “Fauvism” only induces stylization around the mouth as it describes general meaning, like artists “Henri Matisse” and “Kees van Dongen” or “Brutalist painting”. And the same observations when comparing “Chinese Painting” and “Chinese Ink Painting”. In contrast, when a text refers to a specific object or

style, the language ambiguity will disappear. For example, “Lord Voldemort”, “Head of Lord Voldemort”, and “Head of Lord Voldemort in fantasy style” reveals similar stylization results. We also see the similar results concerning the Pixar style. In the interests of brevity, we use “Fauvism” to represent “painting, oil on canvas, Fauvism

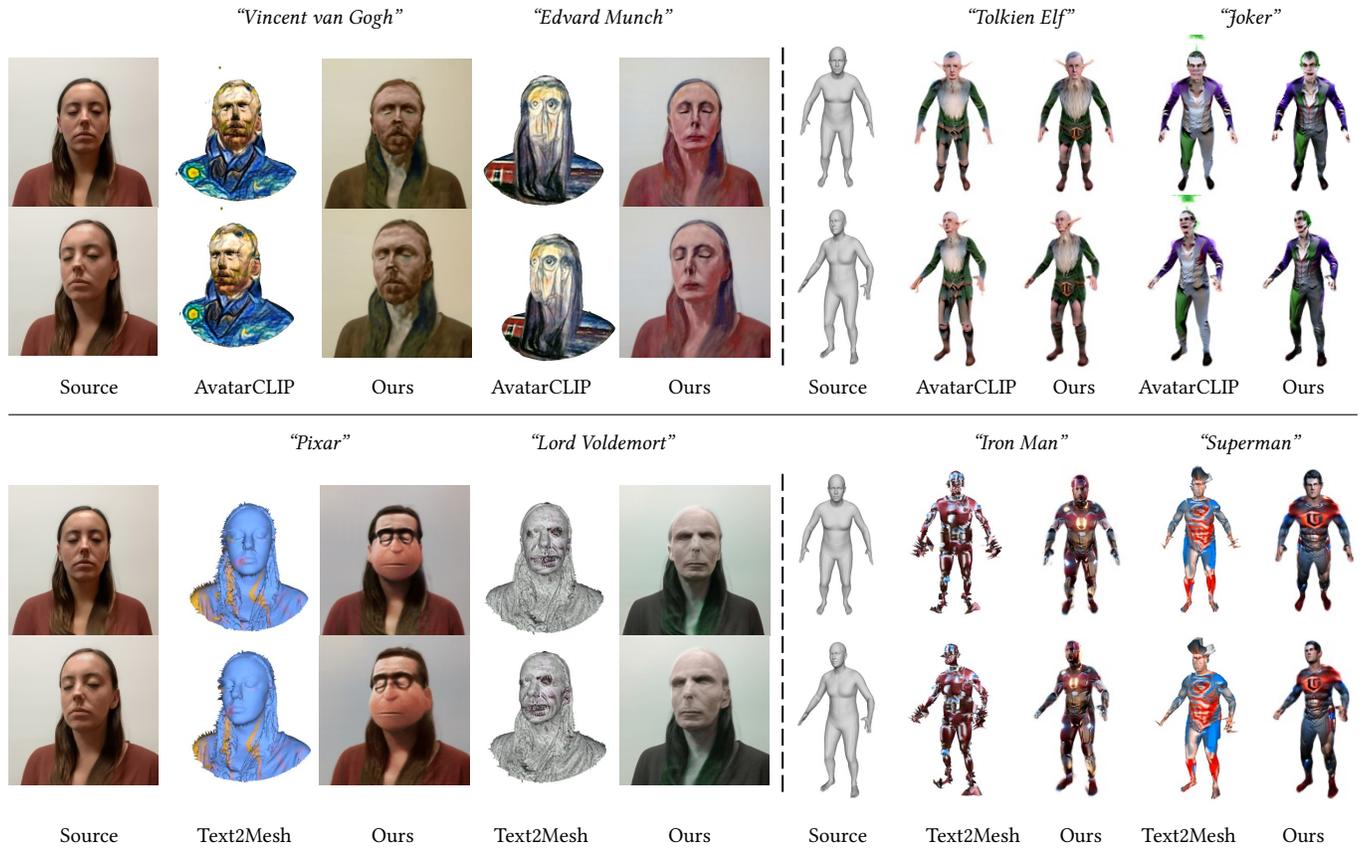


Fig. 7. **Comparisons.** Comparisons to text-guided mesh-based stylization method *Text2Mesh* [Michel et al. 2021] and *AvatarCLIP* [Hong et al. 2022].

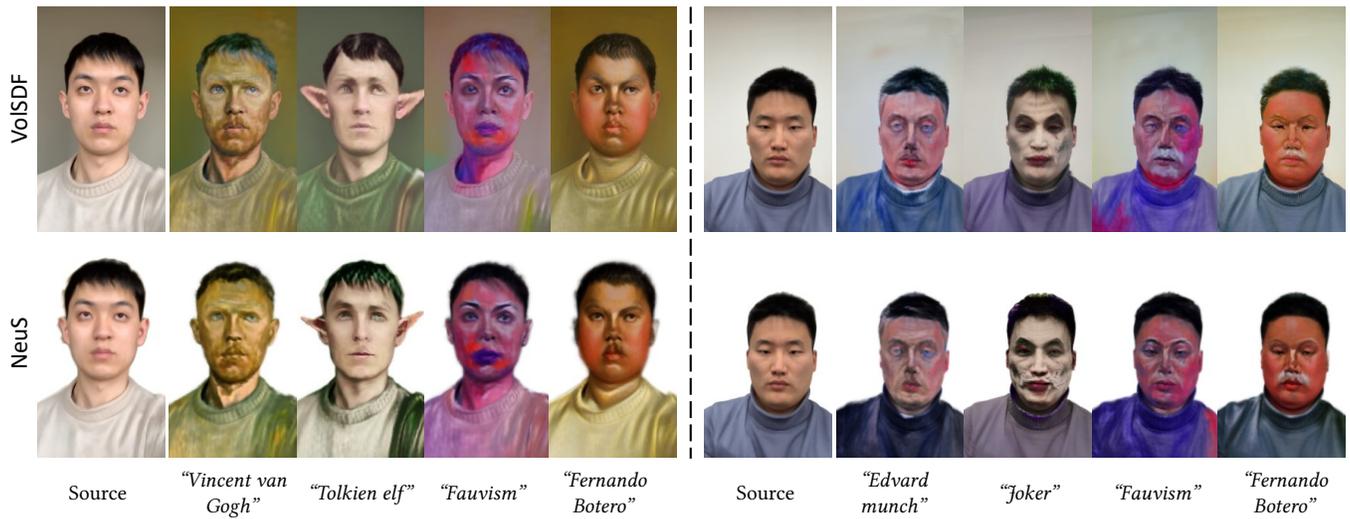


Fig. 8. **Generalization Evaluation.** Generalization evaluation on VoSDF and NeuS.

style" and "Vincent van Gogh" to represent "painting, oil on canvas,

Vincent van Gogh self-portrait style" in other experiments. We also

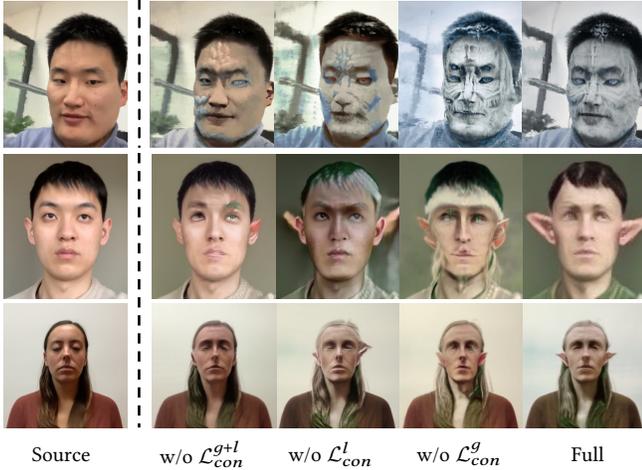


Fig. 9. **Ablations on CLIP-Guided Losses.** Without our global-local contrastive losses, the results suffer from insufficient or non-uniform stylization. The target prompts are “White Walker” and “Tolkien Elf” respectively.

use the same prompt augmentation strategy for other painting styles, including “Edvard Munch” and “Fernando Botero”.

5.4 Comparisons

We compare with most related works following three categories: 1) Text-driven image stylization: *StyleGAN-NADA* [Gal et al. 2021]; 2) Text-driven mesh-based stylization: *Text2Mesh* [Michel et al. 2021] and *AvatarCLIP* [Hong et al. 2022]; and 3) Text-driven NeRF stylization: *CLIP-NeRF* [Wang et al. 2021a] and *DreamField* [Jain et al. 2022]. To make fair comparisons with these methods, we adopt author-released codes and accommodate the input to each method as required. For *StyleGAN-NADA*, we follow its steps to first conduct a face alignment under the setting of *FFHQ* [Karras et al. 2019] and then invert these faces using *e4e* [Tov et al. 2021] into latent codes, before inputting them to *StyleGAN-NADA*. We have also tried *pSp* [Richardson et al. 2021] to invert latent codes but finally adopt *e4e* to obtain better stylization results. Per the authors’ advice, we trained 600 iterations and sampled faces present visual-pleasing stylized results. We place final stylized faces back on the input images by inverting the face alignment process. As for *Text2Mesh*, the input mesh of one example (*‘Lady’*) is provided by the *H3DS* [Ramon et al. 2021], while the input mesh of another example (*‘Human’*) is fetched from *AvatarCLIP*. Both meshes are normalized into -1 to 1, before inputting them to *Text2Mesh*. We follow the training setting of *Text2Mesh* in stylizing the person object to stylize *‘Lady’* and *‘Human’*. We compare to *DreamField* and *AvatarCLIP* following the shape sculpting and texture generation process of *AvatarCLIP*. Similar to *AvatarCLIP*, we also adopt prompt augmentations when stylizing the *‘Human’*. For example, we use text prompts including “Tolkien Elf”, “the back of Tolkien Elf”, and “the face of Tolkien Elf” for the detailed refinement.

The visual comparisons are demonstrated in Figure 5, Figure 6, and Figure 7. For video results, please see the supplementary material.

Comparisons to text-driven image stylization. Compared to *StyleGAN-NADA*, our method can better ensure the desired style strength in all examples by introducing global-local contrastive learning. *StyleGAN-NADA* achieves visual-pleasing results on sampled faces but reflects a degradation for in-the-wild faces partly due to the latent code inversion. Moreover, as a 3D stylization method, ours can preserve view consistencies in the stylized results. In contrast, *StyleGAN-NADA* stylizes each view independently, thus introducing inconsistent shapes or textures to different views. This may lead to flickering artifacts when applied to video applications. Moreover, *StyleGAN-NADA* is less friendly to real faces as the input image has to be inverted back to the *StyleGAN* latent space before stylization, which will inevitably lead to some detail loss and identity change. Unlike it, *NeRF-Art* is not constrained by any latent space of pre-trained networks and does not need the inversion step.

Comparisons to text-driven NeRF stylization. Compared with *CLIP-NeRF*, our advantages are two-fold. First, *CLIP-NeRF* stylizes NeRF using the absolute directional loss, which does not put enough stylizations. Moreover, it suffers from uneven stylizations. For example, we only see enough stylizations on the nose and hair for style “Fauvism”, but the man’s cheek has not been fully stylized. In contrast, we design a global-local contrastive learning strategy to ensure the desired style strength. Second, as no weight regularization is used in *CLIP-NeRF*, its results may appear as severe geometry noises. In contrast, our weight regularization suppresses geometric noises by encouraging a more concentrated density distribution. *DreamField* also adopts the absolute directional loss to stylize NeRF, which cannot guarantee sufficient and uniform stylization. *DreamField* adopts a random background augmentation to CLIP’s attention on the foreground, which requires view-consistent masks, while ours does not. Moreover, our method consistently outperforms *DreamField* in detailed cloth wrinkles, facial attributes, and fine-grained geometry deformations, like muscle shapes and antennas. In summary, our *NeRF-Art* outperforms these methods by proposing a contrastive learning technique to achieve sufficient and uniform stylization and designing a weight regularization to remove cloudy artifacts and geometry noises.

Comparisons to text-driven mesh-based stylization. *Text2Mesh* also supports geometry deformation and texture stylization of a 3D model like ours. However, it assumes there exists a synergy between the input 3D geometry and the target prompt and is more likely to fail when stylizing a 3D mesh towards a less related prompt, such as “Pixar” for the *‘Lady’* model in Figure 6. With carefully-designed loss constraints, ours is more robust to different prompts, either related to the 3D scenes or not. Moreover, limited by the expressivity of the mesh representation, *Text2Mesh* fails most runs and presents unstable stylization results, resulting in irregular deformations and indentations on the edge or surface. Authors of *AvatarCLIP* also report similar results when comparing to *Text2Mesh*. Similar to *DreamField*, *AvatarCLIP* adopts a random background augmentation to lead CLIP to focus on the foreground and prevents floating artifact generations. Nevertheless, this process requires view-consistent masks while ours does not. Moreover, *AvatarCLIP* adds an additional

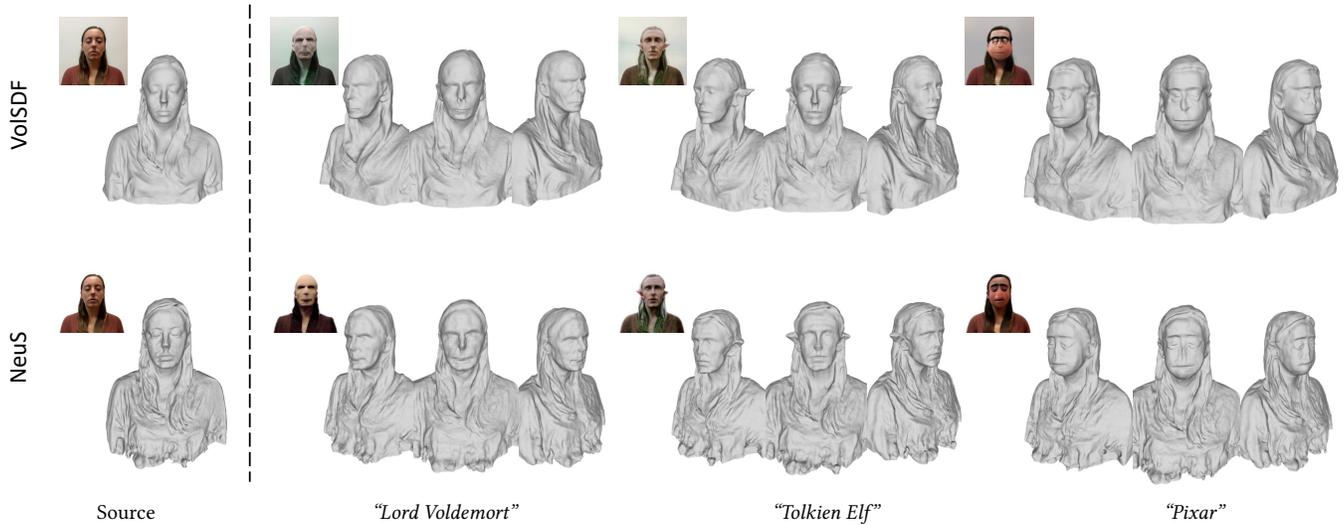


Fig. 10. **Geometry Evaluation.** Our method modulates the geometry and color simultaneously of a pre-trained NeRF to match the desired style described by a text prompt.

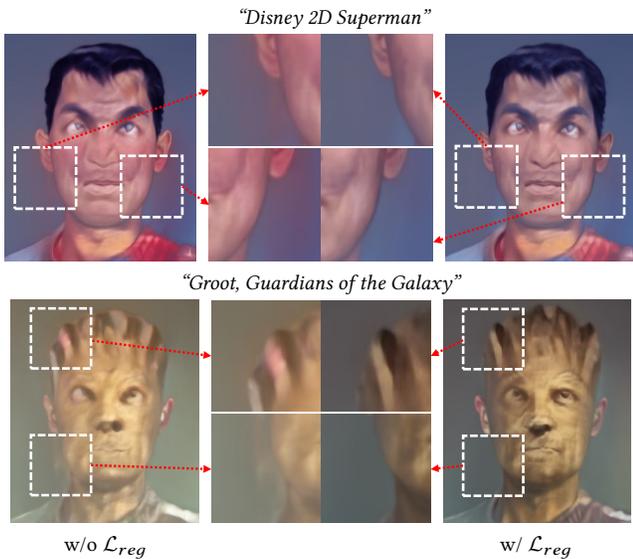


Fig. 11. **Ablations on Weight Regularization.** The cloudy artifacts near the corner or geometric noises are observed without the weight regularization loss.

color network to constrain the general shape of the avatar as well as introducing random shading and lighting augmentations on the textured renderings to strengthen the stylization. Even with these augmentations, *AvatarCLIP* still fails to produce satisfying texture and geometry details. In contrast, ours reveals a fine-grained beard, detailed wrinkles of garments, and clearer face attributes. Noteworthy, our *NeRF-Art* supports stylizing in-the-wild faces, while

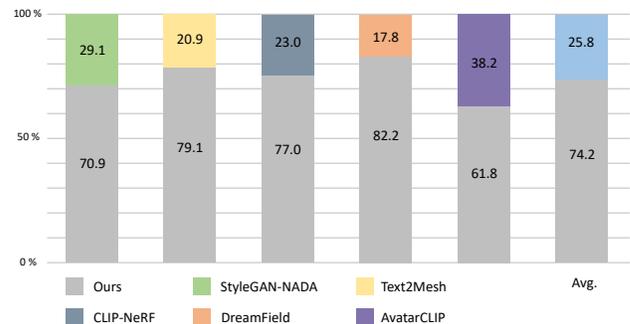


Fig. 12. **User Study.** Our method consistently outperforms state-of-the-art text-guided stylization methods on user preference rates (%).

AvatarCLIP requires a 3D mesh as input to conduct these augmentations. Finally, *AvatarCLIP* can still generate random bumps in the background and make the extracted surface noisy. This is because *AvatarCLIP* sampled a sparse rays (112×112) to construct a coarse renderings for CLIP constraints, due to OOM problem. We found worse results with more noise when reducing sampled ray numbers. In contrast, our method supports training stylization on all rays by imposing a memory-saving technique. In conclusion, *NeRF-Art* achieves better stylization using the proposed contrastive learning strategies without any mesh guidance.

5.5 User Study

To evaluate stylization quality from human perception, we conducted a user study. For each compared category, we used two subjects. For each subject, we selected 5 prompts from our text descriptions dataset and finally obtained 10 test cases for each category and 50 in total. For every test case, we showed one sample of input

frames, the textual prompt, and the results of different methods in two views and random order. The participants were given unlimited time to select the best stylization results by jointly considering three aspects: preservation of the content, faithfulness to the style, and view consistency. We finally collected 23 questionnaires completed by 10 male and 13 Lady participants. Statistics of the user study are shown in Figure 12. Our method outperforms *StyleGAN-NADA*, *CLIP-NeRF*, *Text2Mesh*, *DreamField*, and *AvatarCLIP* by achieving much higher user preference rates.

5.6 Ablation Study

Why global-local contrastive learning? A straightforward way to stylize NeRFs is to apply the directional CLIP loss proposed by StyleGAN-NADA [Gal et al. 2021] to the rendered views. Unfortunately, the directional CLIP loss can enforce the right stylization trajectory but struggles to reach a sufficient magnitude, as shown in the 2nd column of Figure 9. This is because the loss only measures the directional similarity between the normalized embedded vectors but ignores their actual distances. In contrast, our global contrastive loss (3rd column of Figure 9) can ensure the proper stylization magnitude by pushing it as close as possible to the target. However, the global contrastive loss still cannot guarantee a sufficient and uniform stylization of the whole scene. The stylization shows excess on certain parts and insufficiency on others, e.g., insufficient stylized faces and excessively stylized eyes in the *“Tolkien Elf”* example in the 3rd column of Figure 9. This may attribute to the fact that CLIP focuses more attention on regions with distinguishable features than on other regions. Our local contrastive loss helps achieve more balanced stylized results by stylizing every local region of the scene (4th column of Figure 9). However, this local contrastive loss without global information may produce excessive facial attributes, e.g., generating more eyes in the *“White Walker”* example and two left ears in the *“Tolkien Elf”* example. This attributes to insufficient semantics involved in a local patch. This problem can be avoided by adding the global contrastive loss at the same time.

By combining both global and local contrastive loss with the directional CLIP, our method successfully achieves uniform stylization with both correct stylization direction and sufficient magnitude (5th column of Figure 9).

Why weight regularization? Altering the geometry of NeRF may potentially cause cloudy artifacts. In Figure 11, we demonstrate that the weight regularization loss can suppress cloudy artifacts and geometric noises by encouraging a more concentrated density distribution for stylization.

5.7 Generalization Evaluation

We conduct a generalization evaluation on VolSDF and NeuS in Figure 8 to evaluate *NeRF-Art*’s ability in adapting to different NeRF-like models. As NeuS reconstructs a coarse result on our in-the-wild data without a mask may due to inaccurate camera estimations, we conduct a segmentation using RVM [Lin et al. 2022] for better reconstruction and dilate the mask using OpenCV with 3×3 kernel and two iterations to allow geometry variations. In Figure 8, our method presents similar stylization results on VolSDF and NeuS,

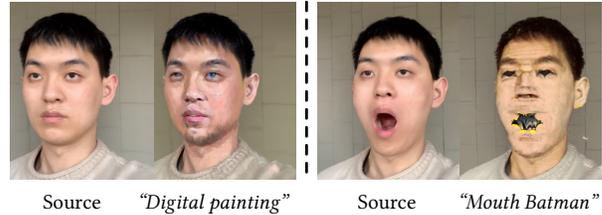


Fig. 13. **Limitations.** Linguistic ambiguity (left) or semantically meaningless words (right) may lead to unexpected results.

which demonstrates that our *NeRF-Art* has the ability to adapt to different NeRF-like models.

5.8 Geometry Evaluation

To evaluate whether the geometry will be correctly modulated in the stylization process, we show the geometry evaluation results in Figure 10. We extract meshes using Marching Cubes [Lorensen and Cline 1987] before and after the stylization for comparison and report results on two widely-used NeRF-like models VolSDF [Yariv et al. 2021] and NeuS [Wang et al. 2021b]. We clearly see geometry changes by comparison with the source mesh. For example, *“Lord Voldemort”* flattens the girl’s nose, *“Tolkien Elf”* sharpens the girl’s ears, and *“Pixar”* rounds the jaw. Moreover, we find the same observations on both VolSDF and NeuS. In summary, we conclude that our method can correctly modulate the geometry of NeRF to match the desired style.

6 CONCLUSION

In this paper, we present *NeRF-Art*, the text-guided NeRF stylization approach based on CLIP. Unlike existing approaches that require the mesh guidance in the stylization process or traps in insufficient geometry deformations and texture details in stylization, ours modulate its geometry and appearance simultaneously to match the desired style and show visual-pleasing results of geometry deformations and texture details with only a text guidance. To achieve it, we introduce a carefully-designed combination of directional constraint to control the style trajectory and novel global-local contrastive loss to enforce the proper style strength. Moreover, we propose a weight regularization strategy to alleviate the cloudy artifacts and geometry noises in deforming the geometry. Extensive experiments on real faces and general scenes show that our method is effective and robust in both stylization quality and view consistency.

Limitations. Despite the success in most cases, our method still has some limitations. First, some text prompts are linguistically ambiguous, like *“Digital painting”*, which describes a wide range of styles, including oil paintings, pencil sketches, 3D rendering images, cartoon drawings, etc. This ambiguity might confuse the CLIP and make the final result unexpected, as shown in Figure 13. Semantically meaningless words cause another kind of unexpected result. For example, if we combine the words *“Mouth”* and *“Batman”* as a prompt, the result unexpectedly puts a bat shape on the mouth, which may not be what the user desires. These are interesting problems worth exploring in the future.

REFERENCES

- Relja Arandjelović and Andrew Zisserman. 2021. Nerf in detail: Learning to sample for view synthesis. *arXiv preprint arXiv:2106.05264* (2021).
- Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. 2021. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5855–5864.
- Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. 2022. Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields. *CVPR* (2022).
- Xu Cao, Weimin Wang, Katashi Nagao, and Ryosuke Nakamura. 2020. Psnnet: A style transfer network for point cloud stylization on geometry and color. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 3337–3345.
- Dongdong Chen, Jing Liao, Lu Yuan, Nenghai Yu, and Gang Hua. 2017a. Coherent online video style transfer. In *Proceedings of the IEEE International Conference on Computer Vision*. 1105–1114.
- Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. 2017b. Stylebank: An explicit representation for neural image style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1897–1906.
- Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. 2018. Stereoscopic neural style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6654–6663.
- Xinghao Chen, Yiman Zhang, Yunhe Wang, Han Shu, Chunjing Xu, and Chang Xu. 2020. Optical flow distillation: Towards efficient and stable video style transfer. In *European Conference on Computer Vision*. Springer, 614–630.
- Pei-Ze Chiang, Meng-Shiun Tsai, Hung-Yu Tseng, Wei-Sheng Lai, and Wei-Chen Chiu. 2022. Stylizing 3D Scene via Implicit Representation and HyperNetwork. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1475–1484.
- Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. 2021. Depth-supervised nerf: Fewer views and faster training for free. *arXiv preprint arXiv:2107.02791* (2021).
- Nianchen Deng, Zhenyi He, Jiannan Ye, Budmonde Duinkharjav, Praneeth Chakravarthula, Xubo Yang, and Qi Sun. 2022. FoV-NeRF: Foveated Neural Radiance Fields for Virtual Reality. *IEEE Transactions on Visualization and Computer Graphics* 28, 11 (2022), 3854–3864.
- Zhiwen Fan, Yifan Jiang, Peihao Wang, Xinyu Gong, Dejia Xu, and Zhangyang Wang. 2022. Unified Implicit Neural Stylization. *arXiv preprint arXiv:2204.01943* (2022).
- Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, and Daniel Cohen-Or. 2021. Stylegan-nada: Clip-guided domain adaptation of image generators. *arXiv preprint arXiv:2108.00946* (2021).
- Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. 2021. Dynamic view synthesis from dynamic monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5712–5721.
- Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. 2021. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 14346–14355.
- Leon A Gatys, Alexander S Ecker, and Matthias Bethge. 2016. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2414–2423.
- Jie Guo, Mengtian Li, Zijing Zong, Yuntao Liu, Jingwu He, Yanwen Guo, and Ling-Qi Yan. 2021. Volumetric appearance stylization with stylizing kernel prediction network. *ACM Trans Graph* 40 (2021), 1–15.
- Fangzhou Han, Shuquan Ye, Mingming He, Menglei Chai, and Jing Liao. 2021a. Exemplar-Based 3D Portrait Stylization. *IEEE Transactions on Visualization and Computer Graphics* (2021). <https://doi.org/10.1109/TVCG.2021.3114308>
- Fangzhou Han, Shuquan Ye, Mingming He, Menglei Chai, and Jing Liao. 2021b. Exemplar-based 3d portrait stylization. *IEEE Transactions on Visualization and Computer Graphics* (2021).
- Aaron Hertzmann. 1998. Painterly rendering with curved brush strokes of multiple sizes. In *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*. 453–460.
- Aaron Hertzmann, Charles E Jacobs, Nuria Oliver, Brian Curless, and David H Salesin. 2001. Image analogies. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*. 327–340.
- Lukas Höllein, Justin Johnson, and Matthias Nießner. 2021. StyleMesh: Style Transfer for Indoor 3D Scene Reconstructions. *arXiv preprint arXiv:2112.01530* (2021).
- Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. 2022. AvatarCLIP: Zero-Shot Text-Driven Generation and Animation of 3D Avatars. In *Proceedings of the ACM SIGGRAPH*.
- Hsin-Ping Huang, Hung-Yu Tseng, Saurabh Saini, Maneesh Singh, and Ming-Hsuan Yang. 2021b. Learning to stylize novel views. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 13869–13878.
- Jialu Huang, Jing Liao, and Sam Kwong. 2021a. Unsupervised image-to-image translation via pre-trained stylegan2 network. *IEEE Transactions on Multimedia* 24 (2021), 1435–1448.
- Xun Huang and Serge Belongie. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*. 1501–1510.
- Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. 2018. Multimodal unsupervised image-to-image translation. In *Proceedings of the European conference on computer vision (ECCV)*. 172–189.
- Yi-Hua Huang, Yue He, Yu-Jie Yuan, Yu-Kun Lai, and Lin Gao. 2022. StylizedNeRF: Consistent 3D Scene Stylization as Stylized NeRF via 2D-3D Mutual Learning. In *Computer Vision and Pattern Recognition (CVPR)*.
- Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. 2022. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 867–876.
- Ajay Jain, Matthew Tancik, and Pieter Abbeel. 2021. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5885–5894.
- Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*. Springer, 694–711.
- Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4401–4410.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8110–8119.
- Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. 2018. Neural 3d mesh renderer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3907–3916.
- Oliver Klehm, Ivo Ihrke, Hans-Peter Seidel, and Elmar Eisemann. 2014. Property and lighting manipulations for static volume stylization using a painting metaphor. *IEEE Transactions on Visualization and Computer Graphics* 20, 7 (2014), 983–995.
- Nicholas Kolkin, Jason Salavon, and Gregory Shakhnarovich. 2019. Style transfer by relaxed optimal transport and self-similarity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10051–10060.
- Hsin-Ying Lee, Hung-Yu Tseng, Qi Mao, Jia-Bin Huang, Yu-Ding Lu, Maneesh Singh, and Ming-Hsuan Yang. 2020. Drit++: Diverse image-to-image translation via disentangled representations. *International Journal of Computer Vision* 128, 10 (2020), 2402–2417.
- Guohao Li, Matthias Müller, Guocheng Qian, Itzel Carolina Delgadillo Perez, Abdullah Abualshour, Ali Kassem Thabet, and Bernard Ghanem. 2021a. Deepgans: Making gans go as deep as cnns. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. 2017. Universal style transfer via feature transforms. *Advances in neural information processing systems* 30 (2017).
- Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. 2021b. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6498–6508.
- Jing Liao, Yuan Yao, Lu Yuan, Gang Hua, and Sing Bing Kang. 2017. Visual attribute transfer through deep image analogy. *arXiv preprint arXiv:1705.01088* (2017).
- Chen-Hsuan Lin, Chen Kong, and Simon Lucey. 2018. Learning efficient point cloud generation for dense 3d object reconstruction. In *proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- Shanchuan Lin, Linjie Yang, Imran Saleemi, and Soumyadip Sengupta. 2022. Robust high-resolution video matting with temporal guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 238–247.
- David B Lindell, Julien NP Martel, and Gordon Wetzstein. 2021. Autoint: Automatic integration for fast neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14556–14565.
- Hsueh-Ti Derek Liu, Michael Tao, and Alec Jacobson. 2018. Paparazzi: surface editing by way of multi-view image processing. *ACM Trans. Graph.* 37, 6 (2018), 221–1.
- Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. 2021. Editing conditional radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5773–5783.
- William E Lorensen and Harvey E Cline. 1987. Marching cubes: A high resolution 3D surface construction algorithm. *ACM siggraph computer graphics* 21, 4 (1987), 163–169.
- Li Ma, Xiaoyu Li, Jing Liao, Qi Zhang, Xuan Wang, Jue Wang, and Pedro V Sander. 2021. Deblur-NeRF: Neural Radiance Fields from Blurry Images. *arXiv preprint arXiv:2111.14292* (2021).
- Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. 2021. Text2Mesh: Text-Driven Neural Stylization for Meshes. *arXiv preprint arXiv:2112.03221* (2021).
- Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. 2019. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 1–14.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields

- for view synthesis. In *European conference on computer vision*. Springer, 405–421.
- Alexander Mordvintsev, Nicola Pezzotti, Ludwig Schubert, and Chris Olah. 2018. Differentiable image parameterizations. *Distill* 3, 7 (2018), e12.
- Fangzhou Mu, Jian Wang, Yicheng Wu, and Yin Li. 2021. 3D Photo Stylization: Learning to Generate Stylized Novel Views from a Single Image. *arXiv preprint arXiv:2112.00169* (2021).
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *arXiv preprint arXiv:2201.05989* (2022).
- Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. 2021. RegNeRF: Regularizing Neural Radiance Fields for View Synthesis from Sparse Inputs. *arXiv preprint arXiv:2112.00724* (2021).
- Michael Niemeyer and Andreas Geiger. 2021. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11453–11464.
- Atsuhiko Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. 2021. Neural articulated radiance field. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5762–5772.
- Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. 2021a. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5865–5874.
- Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. 2021b. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228* (2021).
- Taesung Park, Alexei A Efros, Richard Zhang, and Jun-Yan Zhu. 2020. Contrastive learning for unpaired image-to-image translation. In *European Conference on Computer Vision*. Springer, 319–345.
- Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. 2021. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2085–2094.
- Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. 2021. Animatable neural radiance fields for modeling dynamic human bodies. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 14314–14323.
- Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. 2021. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10318–10327.
- Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. 2017. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 652–660.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020* (2021).
- Eduard Ramon, Gil Triginer, Janna Escur, Albert Pumarola, Jaime Garcia, Xavier Giro-i Nieto, and Francesc Moreno-Noguer. 2021. H3D-Net: Few-Shot High-Fidelity 3D Head Reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5620–5629.
- Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. 2021. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 14335–14345.
- Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. 2021. Encoding in style: a stylegan encoder for image-to-image translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2287–2296.
- Manuel Ruder, Alexey Dosovitskiy, and Thomas Brox. 2016. Artistic style transfer for videos. In *German conference on pattern recognition*. Springer, 26–36.
- Johannes L Schonberger and Jan-Michael Frahm. 2016. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4104–4113.
- Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. 2020. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems* 33 (2020), 20154–20166.
- Bin Sheng, Ping Li, Chenhao Gao, and Kwan-Liu Ma. 2018. Deep neural representation guided face sketch synthesis. *IEEE transactions on visualization and computer graphics* 25, 12 (2018), 3216–3230.
- Yezhi Shu, Ran Yi, Mengfei Xia, Zipeng Ye, Wang Zhao, Yang Chen, Yu-Kun Lai, and Yong-Jin Liu. 2021. Gan-based multi-style photo cartoonization. *IEEE Transactions on Visualization and Computer Graphics* (2021).
- Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. 2021. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7495–7504.
- Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. 2021. Designing an encoder for stylegan image manipulation. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–14.
- Edgar Treitsch, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. 2021. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 12959–12970.
- Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. 2021a. CLIP-NeRF: Text-and-Image Driven Manipulation of Neural Radiance Fields. *arXiv preprint arXiv:2112.05139* (2021).
- Kangan Wang, Sida Peng, Xiaowei Zhou, Jian Yang, and Guofeng Zhang. 2022. NerfCap: Human Performance Capture With Dynamic Neural Radiance Fields. *IEEE Transactions on Visualization and Computer Graphics* (2022).
- Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. 2021b. NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction. *Advances in Neural Information Processing Systems* 34 (2021), 27171–27183.
- Tianyi Wei, Dongdong Chen, Wenbo Zhou, Jing Liao, Zhentao Tan, Lu Yuan, Weiming Zhang, and Nenghai Yu. 2021. Hairclip: Design your hair by text and reference image. *arXiv preprint arXiv:2112.05142* (2021).
- Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. 2021. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9421–9431.
- Guo-Wei Yang, Wen-Yang Zhou, Hao-Yang Peng, Dun Liang, Tai-Jiang Mu, and Shi-Min Hu. 2022. Recursive-NeRF: An efficient and dynamically growing NeRF. *IEEE Transactions on Visualization and Computer Graphics* (2022).
- Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. 2021. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems* 34 (2021), 4805–4815.
- Zipeng Ye, Mengfei Xia, Yanan Sun, Ran Yi, Minjing Yu, Juyong Zhang, Yu-Kun Lai, and Yong-Jin Liu. 2021. 3D-CariGAN: an end-to-end solution to 3D caricature generation from normal face photos. *IEEE Transactions on Visualization and Computer Graphics* (2021).
- Kangxue Yin, Jun Gao, Maria Shugrina, Sameh Khamis, and Sanja Fidler. 2021. 3dstylenet: Creating 3d shapes with geometric and texture style variations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 12456–12465.
- Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. 2021a. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5752–5761.
- Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. 2021b. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4578–4587.
- He Zhang, Fan Li, Jianhui Zhao, Chao Tan, Dongming Shen, Yebin Liu, and Tao Yu. 2022b. Controllable Free Viewpoint Video Reconstruction Based on Neural Radiance Fields and Motion Graphs. *IEEE Transactions on Visualization and Computer Graphics* (2022).
- Kai Zhang, Nick Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snavely. 2022a. ARF: Artistic Radiance Fields. *arXiv preprint arXiv:2206.06360* (2022).
- Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. 2020b. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492* (2020).
- Mohan Zhang, Jing Liao, and Jinhui Yu. 2020a. Deep Exemplar-based Color Transfer for 3D Model. *IEEE Transactions on Visualization and Computer Graphics* (2020).
- Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. 2021. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (TOG)* 40, 6 (2021), 1–18.
- Yandan Zhao, Xiaogang Jin, Yingqing Xu, Hanli Zhao, Meng Ai, and Kun Zhou. 2014. Parallel style-aware image cloning for artworks. *IEEE Transactions on Visualization and Computer Graphics* 21, 2 (2014), 229–240.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*. 2223–2232.