# Geometry Transfer for Stylizing Radiance Fields

Hyunyoung Jung[1*]      Seonghyeon Nam[2]      Nikolaos Sarafianos[2]
Sungjoo Yoo[1]      Alexander Sorkine-Hornung[2]      Rakesh Ranjan[2]

[1]Seoul National University      [2]Meta Reality Labs
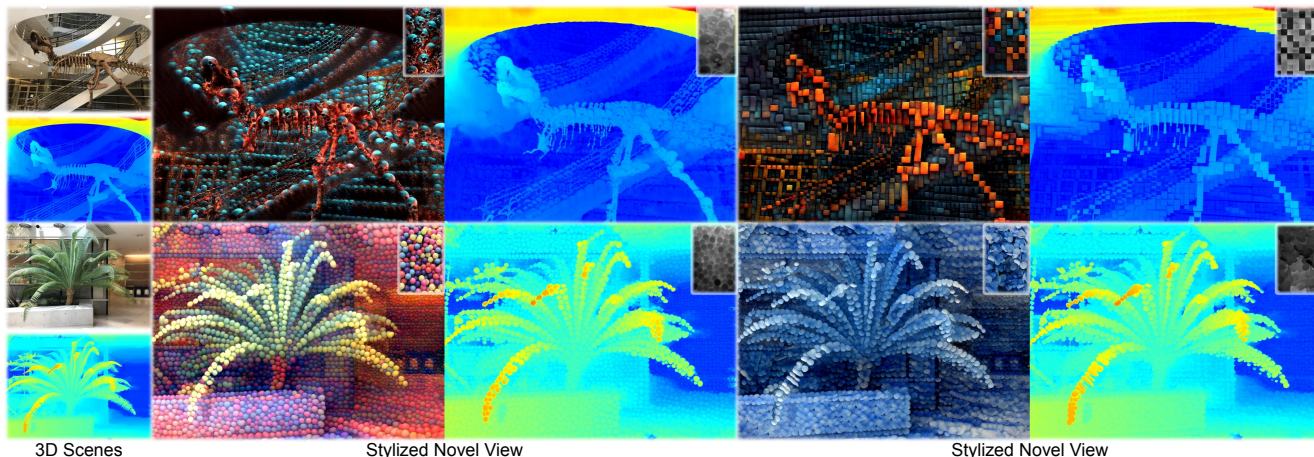
https://hyblue.github.io/geo-srf

Figure 1. Given a reference 3D scene and a pair of style guides: an RGB image and a depth map, we coherently stylize both the scene's appearance and shape to best express the given style.

## Abstract

*Shape and geometric patterns are essential in defining stylistic identity. However, current 3D style transfer methods predominantly focus on transferring colors and textures, often overlooking geometric aspects. In this paper, we introduce Geometry Transfer, a novel method that leverages geometric deformation for 3D style transfer. This technique employs depth maps to extract a style guide, subsequently applied to stylize the geometry of radiance fields. Moreover, we propose new techniques that utilize geometric cues from the 3D scene, thereby enhancing aesthetic expressiveness and more accurately reflecting intended styles. Our extensive experiments show that Geometry Transfer enables a broader and more expressive range of stylizations, thereby significantly expanding the scope of 3D style transfer.*

## 1. Introduction

With the increasing demand for content creation for virtual and augmented reality, style transfer [17] has emerged as an innovative technique that bridges the beauty of art with the precision of technology. At its core, style transfer involves

rendering one image in the stylistic manner of another, producing a new image that combines the foundational structure of the former with the aesthetic qualities of the latter.

In its early phases, style transfer was primarily applied to 2D images [9, 34, 40, 49] and later extended to videos [24, 41, 76, 80] to achieve temporally consistent stylization across image sequences. Recent works have tackled the 3D style transfer problem, by applying styles to 3D models, such as point clouds [25, 54] and meshes [23, 44]. They stand apart from 2D methods, aiming to ensure a cohesive style across multiple camera angles and enabling free-viewpoint rendering. Due to the error-prone geometry stemming from the required 3D reconstruction of them, the stylization of radiance fields [53] has been actively explored. Methods have incorporated global [55] and local [83] constraints, utilized stylized reference views [86], and enhanced diversity through hash encoding and semantic matching [56]. Zero-shot approaches [45] have also been developed to circumvent tedious optimization.

These works focus on transferring aesthetic qualities in terms of colors, texture, and brushstrokes from style images to enhance stylization, effectively applying these attributes to 3D scenes. However, the potential benefits of *geometry* remain largely unexplored and neglected. Even though 3D

---

scenes and objects naturally possess both shape and color attributes, most techniques focus solely on color, leaving the geometric parameters unchanged during the style transfer. Nguyen-Phuoc et al. [55] adjust geometry, but the output shapes do not deviate from the original content and fail to reflect geometric cues from the style images.

As noted by art theorists and image creation experts [1, 22, 33], geometry has historically played a crucial role in defining and influencing style. The shapes and geometric patterns in an artwork are essential to its stylistic identity. From this perspective, in the literature of 2D image style transfer, techniques like correspondence search and image warping [33, 47, 48] have been employed to distort image shapes, showcasing how geometry can enhance the expressiveness of stylization. When applied to 2D images, however, geometric distortion inherently has its limitations. Although shapes are fundamentally 3D forms, images capture only their 2D projections; thus, warping and distorting edges in images offer only an implicit sense of the intended style. It is somewhat limited to assert that they accurately reflect the shape of the objects in the style image.

In this study, we primarily focus on the benefits of incorporating geometric deformation into 3D style transfer. Unlike previous approaches, we define "geometric style" as a distinct and clear characteristic that truly captures the geometric essence of the style image. Our objective is to transfer these intricate forms into the content of a 3D scene. To the best of our knowledge, our work is the first in style transfer literature to propose *Geometry Transfer*, which extracts geometric style from a style template using a depth map and then directly stylizes the shape of neural radiance fields. However, merely replacing the RGB style image from the previous methods with a depth map does not produce the desired results. This issue arises from the intrinsic separation between appearance and shape in the radiance fields representation. When the shape is directly optimized, the resulting colors are not well-aligned with the updated form. To overcome this challenge, we introduce a novel application of deformation fields [61] that predicts the offset vector for each 3D point. This ensures a harmonious stylization of both appearance and shape during optimization. Consequently, we demonstrate the potential of stylizing the geometry of a 3D scene using a depth map as a style image. Beyond this demonstration, we introduce innovative techniques to highlight how stylization can benefit from the incorporation of geometry.

Building on our geometry transfer, we propose a new 3D style transfer method using an RGB-D pair as the style image, aiming for more expressive stylization that better reflects the given style in terms of both shape and appearance. Toward this goal, we propose geometry-aware matching to enhance the diversity of stylization while preserving local geometry through a patch-wise scheme. Additionally, we introduce a novel style augmentation strategy to bring a richer sense of scene depth. Our contributions are summarized as follows:

- For the first time in style transfer literature, we introduce *Geometry Transfer*, a method that extracts style from a depth map and stylizes the geometry of radiance fields.
- We propose a novel usage of deformation fields to ensure coherent stylization of both shape and appearance.
- We introduce novel RGB-D stylization techniques, enhancing expressiveness and better reflecting the style by leveraging scene geometries.
- Our proposed methods can be seamlessly incorporated into existing Panoptic Radiance Fields [68], enabling partial stylization of scenes for more practical applicability.

## 2. Related Works

**Neural style transfer.** Neural style transfer is the process of creating a new image that fuses the structural elements of a content image with the aesthetic characteristics of a style image. Gatys et al. [17] described the style transfer as an iterative optimization that aligns feature correlations from both images, using a deep feature network [69]. Building on this, various techniques [5, 9, 34, 36–38, 40, 49] have advanced stylization through semantic correspondence [27, 43, 85], image blending [4, 32, 50, 71, 84] and novel loss formulas for feature statistics computation [21, 51, 62]. To address the slow convergence of iterative optimization, a feed-forward network has been widely adopted to facilitate arbitrary stylization [12, 19, 26, 42, 46, 58, 67] in real-time [29, 64, 77, 79, 80]. In response to inconsistent stylization across multiple images, techniques to stylize videos [7, 15, 24, 41, 63, 75, 76, 80] and stereo images [8, 18] have been proposed. Given that shape and geometry are essential for expressive stylization, certain methods have focused on distorting the content's structure, specifically for faces [82] and text [81]. More general methods [33, 47, 48] utilize correspondence searches and image warping to align content and style from the same class identity. Since shapes and geometry are fundamentally 3D forms, however, modifications to 2D images often fail to capture the accurate style of geometry. Our proposed method aims to directly stylize the shape of the 3D scene using the estimated depth map from the style image.

**3D style transfer.** Recent techniques have applied style transfer to 3D models to ensure coherent stylization across images rendered from multiple viewpoints. Earlier methods stylized explicit representations, such as point clouds [25, 54], and mesh [23, 44]. More recent techniques [10, 14, 45, 55, 56, 83, 86, 87] have actively explored the stylization of implicit representations, i.e. radiance fields [53]. In optimization-based approaches [14, 56], Nguyen-Phuoc et al. [55] alternated between rendering and 2D stylization, using a global style loss to stylize the 3D scene. Meanwhile,
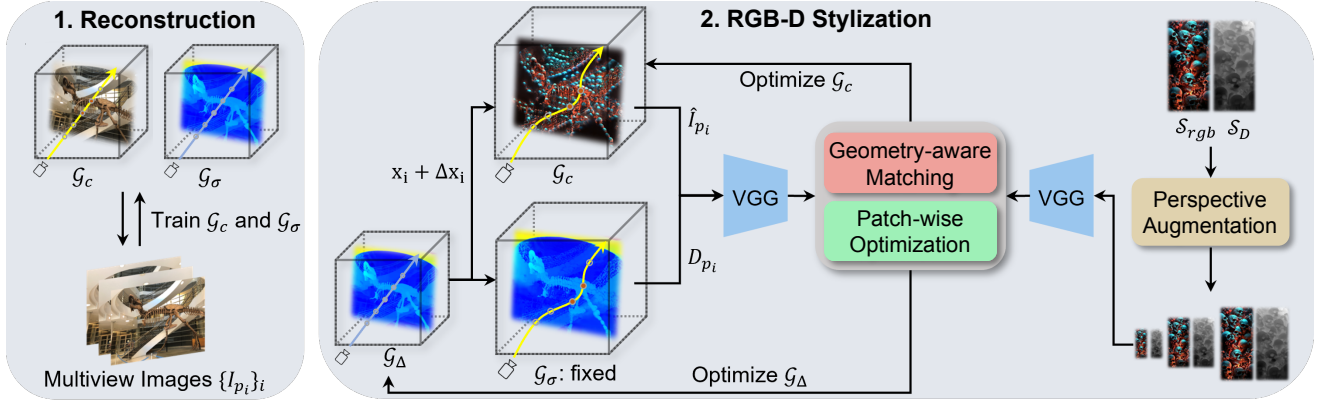
Figure 2. **Overview of our method.** First we pre-train TensoRF [6] on real-world images to obtain the color grid $\mathcal{G}_c$ and density grid $\mathcal{G}_\sigma$, enabling photorealistic reconstruction. Subsequently, we extract VGG features from style images as an RGB-D pair to stylize the shape and appearance of radiance fields. Here, the shape is modified through the additional deformation grid $\mathcal{G}_\Delta$, while $\mathcal{G}_\sigma$ remains fixed.

Zhang et al. [83, 86] employed a nearest-neighbor matching loss [35] and utilized a reference stylized view [86] to enhance detail preservation. Another direction avoids per-style optimization and instead adopts hypernetworks [10], feature transformations [45], and Lipschitz mappings [87] to facilitate arbitrary stylization of 3D scenes. While most techniques prioritize appearance without altering geometry during style transfer, we emphasize geometry distortion to improve stylization expressiveness and style accuracy. To our knowledge, this is the first use of a depth map as a style guide to optimize the radiance fields' geometry. Instead of using reference images for style as above, several approaches stylize radiance fields using text prompts via CLIP embedding [72, 73] and leverage diffusion models [20, 31].

**Deformation fields.** Deformation fields have been widely used initially to model the target 3D shape of objects while preserving their geometric details [13, 30, 74]. They define the shape as a surface deformation of the template 3D models. Pumarola et al. [61] introduced D-NeRF, which employs a time-varying deformation function to capture the transformation between canonical and deformed scenes. This approach allows for the reconstruction of dynamic scenes using a single moving camera. Building on this concept, subsequent studies [16, 59, 60, 70] have addressed the view synthesis challenge in dynamic scenes. Our approach distinguishes itself by using deformations to ensure alignment of shape and appearance when stylizing the geometry of the radiance fields.

## 3. Methodology

**Preliminaries: Stylizing Radiance Fields.** Stylizing radiance fields is conceptualized as a fine-tuning process that begins with a pre-trained NeRF on a real-world 3D scene. We use TensoRF [6] as our scene representation. It introduces two separate grids, $\mathcal{G}_c$ and $\mathcal{G}_\sigma$, each with per-voxel multi-channel features where the former models appear-

ance, and the latter the volume density. To ensure efficient rendering and compact representation, TensoRF factorizes them into multiple low-rank components. For pre-training on the target 3D scene, which includes training images $\{I_i\}_{i=1}^N$ and their corresponding camera poses $\{p_i\}_{i=1}^N$, we follow the training scheme outlined in the original paper and refer the reader there for additional details.

We primarily follow the methods of stylizing radiance fields in ARF [83]. In each stylization iteration, we randomly select a viewpoint $p_i$ and render the image $\hat{I}_{p_i}$. We then extract 2D feature maps $F_\mathcal{I}^{rgb}$ from $\hat{I}_{p_i}$ and $F_\mathcal{S}^{rgb}$ from the style image, $\mathcal{S}_{rgb}$, using VGG [69]. After this, we compute the style loss $L_{style}$ between these feature maps, formulated as a nearest-neighbor matching loss [35, 83]:

$$L_{style} = \frac{1}{N} \sum_{i,j} \min_{i',j'} D(F_\mathcal{I}^{rgb}(i,j), F_\mathcal{S}^{rgb}(i',j')), \quad (1)$$

where $D(,)$ computes the cosine distance between two normalized feature vectors.

### 3.1. Geometry Transfer

Our approach stems from a fundamental question: Can we transfer *geometry* in the same manner as we transfer colors? To explore this, we introduce the use of a depth map as a style guide to transfer its geometry to a 3D scene. An overview of our approach is depicted in Fig. 2.

#### 3.1.1 Depth Map as a Style Guide

Instead of using an RGB image as the style guide, we replace it with a depth map, denoted as $\mathcal{S}_D$, which captures a distinct style of shape. During the style transfer process, we render the depth map $D_{p_i}$ and optimize the style loss between $D_{p_i}$ and $\mathcal{S}_D$. Since the VGG network expects 3-channel images as input, we concatenate the depth maps
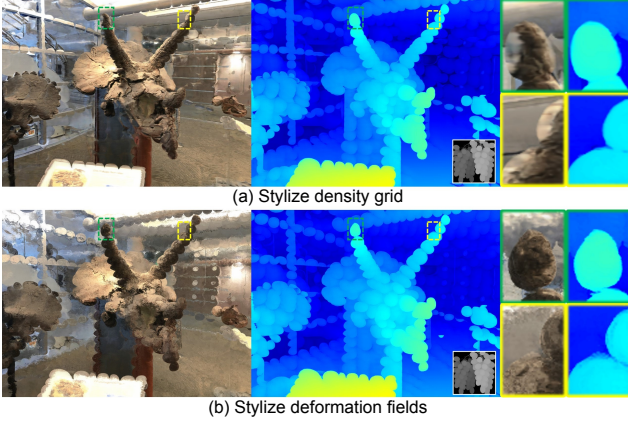
Figure 3. **Comparisons of the stylized results** obtained by optimizing the density grid (a), and by optimizing the deformation fields (b). When directly optimizing the density, background colors are assigned to the updated parts of the foreground object.

along the channel dimensions by replicating them three times. Given that $D_{p_i}$ relates solely to volume density, the loss function optimizes the density grid, $\mathcal{G}_\sigma$. As illustrated in Fig. 3 (a), this approach revealed that we could manipulate shapes in the same manner that we apply style transfer to colors. However, a challenge arises: while the shape adapts to the style image, the color fields remain static, leading to undesired outcomes. For instance, background colors might be applied to updated portions of foreground objects, even though ideally, the colors of these objects should evolve cohesively with their shape.

#### 3.1.2 Modeling Deformation Fields

After pre-training on real-world scenes, the density grid $\mathcal{G}_\sigma$ forms a surface distribution that mirrors the target 3D scene. Concurrently, color values in the appearance grid $\mathcal{G}_c$ are updated in coherence with the corresponding locations of the surface distribution in $\mathcal{G}_\sigma$. This synchronization leads to the rendering of precise surfaces with accurate appearance. However, when the geometry is stylized, the surface distribution within $\mathcal{G}_\sigma$ changes, yet $\mathcal{G}_c$ remains consistent. During sampling of 3D points along rays and querying colors and densities from these misaligned grids, the colors of the modified areas are predominantly sourced from the new surface locations in $\mathcal{G}_c$, as shown in Fig. 4 (a), even though the color fields still align with the original distribution.

To address this issue, we introduce a deformation network to enable synchronous modifications of both shape and appearance. This network is designed as a function predicting a three-dimensional displacement vector, $\Delta\mathbf{x}_i \in \mathbb{R}^3$, that maps a 3D point $\mathbf{x}_i$ to its canonical location $\mathbf{x}_i + \Delta\mathbf{x}_i$. In our context, the canonical space refers to the original scene before stylization. We represent the deformation network using another voxel grid, $\mathcal{G}_\Delta$, and update it ex-
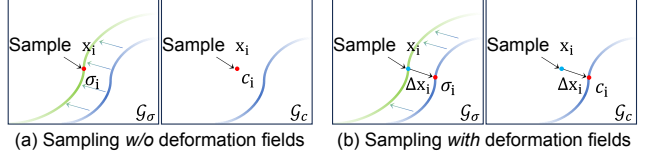


(a) Sampling *w/o* deformation fields    (b) Sampling *with* deformation fields

Figure 4. **Sampling w/ and w/o deformation fields.** Comparisons of the sampling density $\sigma_i$ and color $c_i$ for a 3D point $\mathbf{x}_i$ with and without deformation fields. The curves represent the 2D projected surfaces of objects, where green depicts the stylized surface and blue the original surface. By sampling with deformation fields, we coherently sample both values from the original surface.

clusively for the purpose of stylizing the geometry, keeping $\mathcal{G}_\sigma$ unchanged. After the stylization, the canonical surface remains intact. When rendering the stylized scene, both the densities and colors are sampled from the original surface locations, as described in Fig. 4 (b). This ensures that coherent colors are associated with the modified areas, leading to the differences shown in Fig. 3 (b).

### 3.2. RGB-D Stylization

To realize a more expressive stylization that modifies both colors and geometry, we employ a pair of style guides: an RGB image and a depth map. Given an RGB style $\mathcal{S}_{rgb}$, we use a zero-shot depth estimation network [3] to derive its depth map. This then serves as the style depth, $\mathcal{S}_\mathcal{D}$.

#### 3.2.1 Geometry-aware Nearest Matching

To stylize using two style images, specifically $\mathcal{S}_{rgb}$ and $\mathcal{S}_D$, the style loss must be adjusted to account for multiple style sources. Since our goal is to align both colors and geometry, computing the nearest matching loss independently is inappropriate due to potential inconsistencies between patterns of appearance and shape. A more effective method is to initially identify the closest match between content and style features in one domain, then compute the style loss for the other domain using these predetermined pairs. Alternatively, both color and geometry features could be used to search for the nearest neighbors concurrently. After extracting VGG feature maps from the two modalities, we concatenate them along the channel dimension and then perform a search to find the nearest pair:

$$j = \arg\min_{i'} D([F_\mathcal{I}^{rgb}(i), F_\mathcal{I}^D(i)], [F_\mathcal{S}^{rgb}(i'), F_\mathcal{S}^D(i')]) \quad (2)$$

We then optimize the cosine distance $D$ separately for features from each modality:

$$L(i) = D(F_\mathcal{I}^{rgb}(i), F_\mathcal{S}^{rgb}(j)) + D(F_\mathcal{I}^D(i), F_\mathcal{S}^D(j)), \quad (3)$$

The style loss is calculated as the mean across all feature vectors: $L_{style} = \frac{1}{N}\sum_i L(i)$. This strategy, which involves incorporating geometry features into the matching process, not only enhances diversity but also better preserves scene structure, as demonstrated in Sec. 4.

4

### 3.2.2 Patch-wise Optimization

With an RGB style image, it is straightforward to determine if the output aligns with the style in terms of color, texture, and other visual attributes. However, in geometry, depth maps provide limited cues to identify the style. This is because shapes are defined not by isolated pixels, but by their relationship to their surroundings. The existing nearest matching loss, which conducts matching on a per-pixel basis, is not enough for transferring the style of geometry effectively. To address this, we introduce a patch-wise matching scheme that broadens the receptive fields, thereby becoming more effective in capturing spatial interactions.

Given the extracted VGG feature maps $F_\mathcal{I}$ and $F_\mathcal{S}$, we first partition each feature map into sets of $k \times k$ patches: $\{\mathcal{P}_\mathcal{I}^i\}_i$ and $\{\mathcal{P}_\mathcal{S}^i\}_i$. The patch-wise style loss $L_{\mathcal{SP}}$ is then given by:

$$L_{\mathcal{SP}} = \frac{1}{|\mathcal{P}_\mathcal{I}|} \sum_i \min_j D^\mathcal{P}(\mathcal{P}_\mathcal{I}^i, \mathcal{P}_\mathcal{S}^j), \qquad (4)$$

where $D^\mathcal{P}(\mathcal{P}_1, \mathcal{P}_2)$ computes the sum of the cosine distances between feature vectors at corresponding locations within each patch:

$$D^\mathcal{P}(\mathcal{P}_1, \mathcal{P}_2) = \sum_i^{k^2} D(F_1^i, F_2^i), \qquad (5)$$

Here, $D$ calculates the cosine distance, and $F_{1,2}$ represents feature vectors that constitute each patch. To achieve larger receptive fields without increasing computation, each patch can be defined with a dilation rate $r$ as a hyperparameter.

### 3.2.3 Perspective Style Augmentation

We typically select style images with distinct patterns as shown in Fig. 1, since this aids in the clearer identification of their geometric style. To enhance diversity and the perception of depth, we can vary the sizes of these patterns, applying them differently to surfaces based on their distance.

Before the stylization process, we gather 3D points in world coordinates from all training viewpoints and categorize them into $N$ bins, $\{B_i\}_{i=1}^N$, based on their $z$-coordinates. Each bin $B_i$ is linked to a central value $C_i$, determined by averaging the $z$ values of points within that bin. Given that pattern sizes can vary with the relative resolutions of content and style images [28], we modify the style images by downsampling them at multiple scales $\{s_i\}_{i=1}^N$. This process results in a series of style pairs, $\{\mathcal{S}_i\}_{i=1}^N$, where $\mathcal{S} = (\mathcal{S}_{rgb}, \mathcal{S}_D)$. We set the scale of the first bin, $s_1$, to 1. To reflect real-world conditions, the scales of subsequent bins are calculated based on their relative distance from the first bin as: $s_i = C_1/C_i$.

During stylization, each pixel in the rendered image is assigned to a bin $B_{i'}$, based on the shortest distance from the pixel's $z$-coordinate to the center $C_{i'}$ of the bin. This method transforms the rendered image into a format akin to a layered depth image [65]. Each layer is then stylized using its corresponding style pair $\mathcal{S}_{i'}$, which is downsampled to the appropriate scale. Consequently, larger patterns are mapped onto surfaces closer to the viewer, while smaller patterns are applied to more distant surfaces, thereby enhancing the overall sense of depth.

## 4. Experiments

**Implementation Details.** We implemented our work based on the code of ARF [83], using TensoRF [6] as the underlying NeRF representation. For the training of TensoRF during the photorealistic reconstruction stage, we followed the training scheme from its original paper and utilized a distortion regularizer [2] to mitigate artifacts such as floaters and background collapse. In this stage, the deformation fields are randomly initialized and are optimized to output zeros for all sampled input points. After pre-training, we maintained the density grid at a constant but updated both the appearance and deformation grids to stylize the reconstructed 3D scene using style loss. We employed the `conv2` and `conv3` layers of the VGG-16 [69], when computing the style loss. We applied the view-consistent color transfer [83] before and after the stylization.

**Datasets.** We conducted experiments on the LLFF dataset [52], comprising high-resolution captures of real-world, forward-facing scenes, as used in recent 3D style transfer methods [55, 83, 86]. Furthermore, we utilized the ScanNet dataset [11] to verify the capability of our approach on scenes captured from diverse camera viewpoints, highlighting our method's potential for partial stylization. The ScanNet dataset includes multiple sequences of real-world indoor scenes characterized by varied trajectories and a collection of common furniture types.

**Evaluation metric.** We use Single Image Fréchet Inception Distance (SIFID) [66] to evaluate the stylizations. SIFID calculates the feature distance between two images, indicating the style similarity between the style image and the stylized results for quantitative evaluation in image style transfer [57, 78]. We introduce three methods to assess how both the shape and appearance reflect the specified style guide:

- *RGB*: To evaluate stylization, we compute the SIFID between the RGB style image, $\mathcal{S}_{rgb}$, and the rendered RGB image, $\hat{I}$.
- *Gray*: Recognizing that shape and pattern forms, beyond color, influence style, we convert both $\mathcal{S}_{rgb}$ and $\hat{I}$ to grayscale. We then compute the SIFID between these images, allowing us to exclude the influence of color and measure the other style elements.
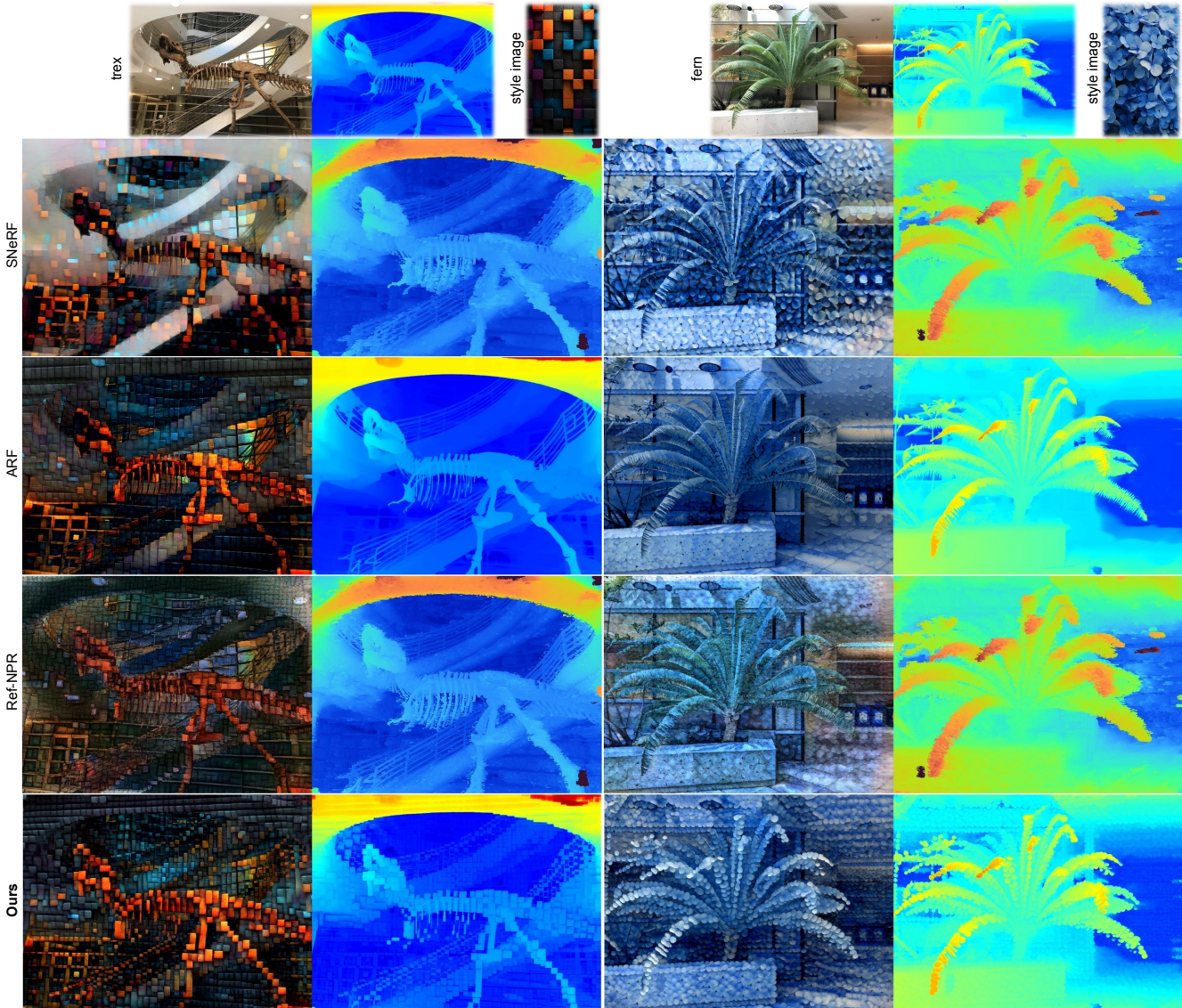- *Depth*: To evaluate the geometry style, we compute the

Figure 5. **Qualitative comparisons** with SNeRF [55], ARF [83] and Ref-NPR [86] on the `trex` and `fern` scenes [52].

SIFID between the depth style $\mathcal{S}_D$ and the rendered depth map, $\hat{D}$.

## 4.1. Qualitative and Quantitative Comparisons

In Fig. 5, we qualitatively compare our results with recent top-performing 3D style transfer methods, including SNeRF [55], ARF [83], and Ref-NPR [86] on the the `trex` and `fern` scenes [52]. All these methods stylize radiance fields, guided by a single style image. The scale of the style image plays a crucial role in replicating the patterns from the style image; hence, we manually tuned these methods to find the optimal configurations. Since SNeRF did not provide an official implementation, we used an alternative version provided by Zhang *et al.* [86], enabling a density update as mentioned in their original paper. For Ref-NPR, we

utilized NNST [35] to generate a reference stylized view. For ARF, we used the authors' provided TensoRF version since the geometry of the scene is noisy and very inaccurate in the original version with Plenoxels. We applied the distortion regularizer [2] to refine its geometry during pretraining for fair comparisons.

As shown in the figure, our method provides clearer colors and more accurately stylized shapes. Notably, our stylized results replicate the clear and complete forms of style patterns, an ability not achievable by merely stylizing colors, due to the limited space available for mapping complete patterns without altering geometry. To be specific, in order to stylize the `fern` leaves, it is necessary to change their shape because the leaves are sharp and narrow, which cannot display patterns of the style image without a shape

| Method | trex | | | fern | | |
|---|---|---|---|---|---|---|
| | RGB | Gray | Depth | RGB | Gray | Depth |
| SNeRF [55] | 1.62 | 0.81 | 0.59 | 1.32 | 0.64 | 0.40 |
| ARF [83] | 1.54 | 0.64 | 0.51 | 1.11 | 0.48 | 0.36 |
| Ref-NPR [86] | 1.59 | 0.72 | 0.61 | 1.75 | 0.79 | 0.41 |
| Ours | **1.43** | **0.58** | **0.44** | **0.81** | **0.37** | **0.28** |

Table 1. **Quantitative comparisons** of SIFID [66] for RGBs, grayscale images, and depth maps with recent methods. Lower scores indicate better performance. For each scene, images are rendered from 30 viewpoints, and their average score is computed.

deformation. Our method accurately stylizes those regions while the others are limited to stylizing only appearance to just hallucinate the shape. Even though SNeRF updates the density during stylization, the resulting geometry does not reflect any cues from the style image because it lacks proper guidance for geometry style.

In Table 1, we compare the SIFID [66] to measure the style similarity between the style images and the rendered images. Our method outperforms others in all metrics, encompassing both appearance and geometry. This demonstrates that incorporating stylization into geometry, as well as colors, enhances the overall style representation and more accurately reflects the intended styles.

In Table 2, we present the results of a user study designed to assess visual appeal based on user preferences. We collected rankings from 22 participants for each set of stylization results produced by Ref-NPR [86], ARF [83], and SNeRF [55], and then computed the average rankings for 12 different stylized scenes. Notably, our proposed method outperforms the others, achieving the highest average ranking. Furthermore, out of 264 total responses ($22 \times 12$), our method was mostly favored, being selected as the best in 162 instances.

| metric | Ours | Ref-NPR | ARF | SNeRF |
|---|---|---|---|---|
| Avg. rank $\downarrow$ | **1.55** | 3.17 | 2.58 | 2.70 |

Table 2. **User study results** reporting the average ranking.

## 4.2. Ablation Experiments

**Geometry-aware nearest matching.** In Fig. 6, we compare the results of the nearest matching exclusively with the color features extracted from $\mathcal{S}_{rgb}$ and our proposed geometry-aware strategy. When using only color features, the resulting style includes similar colors and overlapping patterns in regions with a similar appearance. This approach tends to wash out object boundaries, rendering them indistinguishable, and leads to a loss of content structure and diversity, particularly in semi-transparent objects. In contrast,



3D Scene     *w/o* geometry feature     *with* geometry feature
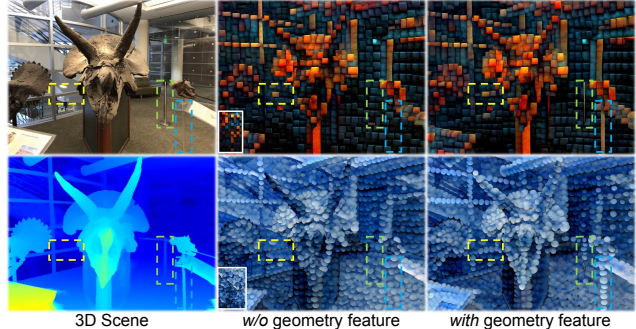
Figure 6. **Ablation study on the impact of geometric features.** Comparison of the results using nearest matching based on color features versus geometry-aware matching. Geometry features enhance diversity and enable distinct stylizations, differentiating objects with similar colors.
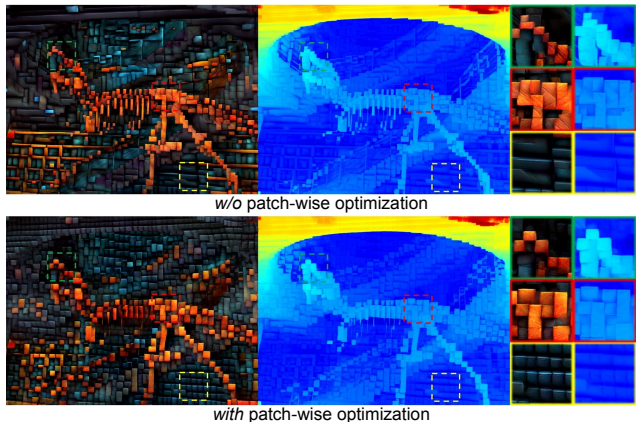


*w/o* patch-wise optimization

*with* patch-wise optimization

Figure 7. **Impact of patch-wise optimization**. The patch-wise scheme enhances the clarity and accuracy of patterns and shapes.

when geometry features are also used, the combined consideration of shape and color during the matching process results in distinct colorization and patterns across objects, even those with similar appearances. This differentiation of boundaries enhances content preservation.

**Patch-wise optimization.** In Fig. 7, we compare the results of the nearest neighbor loss with and without our proposed patch-wise optimization. Without the patch-wise scheme (the top figure), each feature vector in the content and style feature maps is independently matched based on the respective cosine distances, which leads to a failure in maintaining local geometry. Due to its small receptive fields, the scene often contains only incomplete parts of patterns and shapes, resulting in decreased style accuracy. In contrast, when applying the patch-wise optimization (the bottom figure), the positions of local neighbors within the feature maps are preserved during the matching process, enabling the capture of larger receptive fields. This approach results in the intact and complete reproduction of patterns from the style image.

**Perspective style augmentation.** In Fig. 8, we compare the effects of our proposed perspective style augmentation on
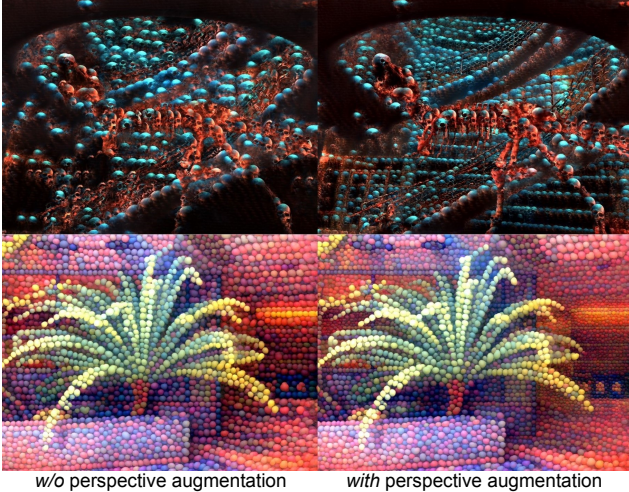
Figure 8. **Perspective style augmentation impact**. The proposed augmentation enhances depth perception by mapping larger patterns to closer surfaces and smaller patterns to more distant ones.
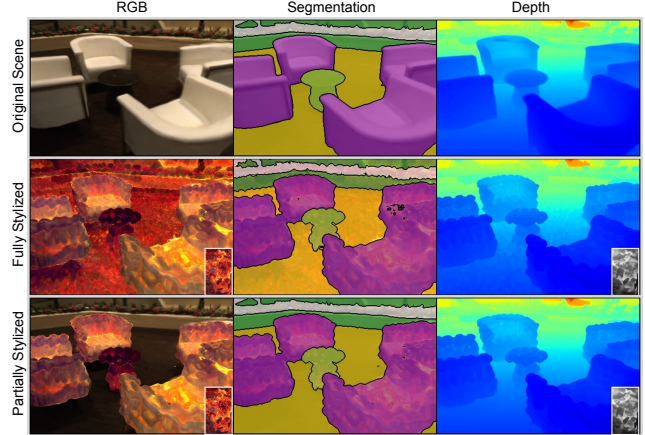
Figure 9. **Stylization with 3D semantic lifting.** We stylize Panoptic Lifting, pre-trained on the ScanNet dataset [11], allowing for the free alteration of target objects for stylization during runtime. As the stylization alters the colors and shapes of objects, the segmentation adapts to their updated forms.

stylization. As depicted on the left, stylizing the entire surface with patterns of uniform size, regardless of the distance from the camera, violates perspective rules and diminishes the sense of depth. Conversely, the right column figures demonstrate that decreasing the pattern sizes based on their depth location enhances the perception of depth in the 2D rendered image and aids in preserving the scene's structure.

### 4.3. Application: Partial Stylization

In Fig. 9, we demonstrate the applicability of our method in partially stylizing 3D scenes. Instead of partially optimizing the scene based on semantic masks [39], we have integrated our method with Panoptic Lifting [68]. This approach involves volumetric representations that produce view-consistent panoptic segmentations, along with RGB values and shapes. Our method can be seamlessly incorporated into it, enabling us to dynamically render and select target classes and object instances during *runtime*.

The Panoptic Lifting models a function that maps a 3D point $\mathbf{x}_i$ to color $c_i$, volume density $\sigma_i$, semantic class probability $\kappa_i$, and object id distribution $\pi_i$ over each class as: $\kappa_i(k)\pi_i(j)$. The underlying representation adopts TensoRF and consists of a color grid $\mathcal{G}_c$ and $\mathcal{G}_\sigma$. To stylize Panoptic Lifting, we introduce an additional deformation grid $\mathcal{G}_\Delta$ and apply our proposed RGB-D stylization methods to optimize both $\mathcal{G}_c$ and $\mathcal{G}_\Delta$. After the style transfer, we obtain the stylized color grid $\mathcal{G}_{c'}$ and use it to freely render the stylized scene. It is important to note that even if the shape changes after stylization, our use of deformation fields enables sampling from the canonical space for color and density, as well as for the classes and object ids. Thus, if the stylization alters the original shapes, the semantic predictions cohesively adapt to the new shapes.

To render the partially stylized view for specific target classes or objects, we begin by estimating the target class for each 3D point along the rays. During RGB rendering, color is sampled for the 3D points that comprise the target objects from the stylized grid $\mathcal{G}_{c'}$. This sampling is conducted after applying deformation to the points, denoted as $\mathbf{x}_i + \Delta\mathbf{x}_i$. The rest of the scene is rendered using colors from the original grid $\mathcal{G}_c$, with no deformation applied.

**Limitations and future work.** Our selection of TensoRF [6] as the underlying representation inherently constrains our capabilities in handling $360°$ unbounded scenes. Also, additional challenges arise due to our focus on accurately transferring the shapes and patterns from a single style image to the 3D scene. This task is highly ill-posed as the patterns in 3D scenes do not appear identical when viewed from significantly different perspectives. To effectively stylize $360°$ scenes, it would be beneficial to investigate the use of multi-view style guides or 3D style guides, extending beyond a single-image style reference.

## 5. Conclusion

We proposed Geometry Transfer, a novel method that uses a depth map as a stylistic guide for modifying the geometry of radiance fields. By innovatively employing deformation fields, we achieved coherent alteration of both shape and appearance in 3D scenes. Building upon this foundation, we developed novel RGB-D stylization techniques, leveraging geometric cues to enhance aesthetic expressiveness and more accurately reflect intended styles. Extensive experiments have shown that our methods facilitate a broader spectrum of stylizations compared to previous approaches, significantly expanding the scope of 3D style transfer.

# References

[1] Rudolf Arnheim. Art and visual perception: A psychology of the creative eye. *Univ of California Press*, 1954. 2

[2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, 2022. 5, 6

[3] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 4

[4] Junyan Cao, Yan Hong, and Li Niu. Painterly image harmonization in dual domains. In *AAAI*, 2023. 2

[5] Bindita Chaudhuri, Nikolaos Sarafianos, Linda Shapiro, and Tony Tung. Semi-supervised synthesis of high-resolution editable textures for 3d humans. In *CVPR*, 2021. 2

[6] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *ECCV*. Springer, 2022. 3, 5, 8

[7] Dongdong Chen, Jing Liao, Lu Yuan, Nenghai Yu, and Gang Hua. Coherent online video style transfer. In *ICCV*, 2017. 2

[8] Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. Stereoscopic neural style transfer. In *CVPR*, 2018. 2

[9] Tian Qi Chen and Mark Schmidt. Fast patch-based style transfer of arbitrary style. *arXiv preprint arXiv:1612.04337*, 2016. 1, 2

[10] Pei-Ze Chiang, Meng-Shiun Tsai, Hung-Yu Tseng, Wei-Sheng Lai, and Wei-Chen Chiu. Stylizing 3d scene via implicit representation and hypernetwork. In *WACV*, 2022. 2, 3

[11] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017. 5, 8

[12] Yingying Deng, Fan Tang, Weiming Dong, Wen Sun, Feiyue Huang, and Changsheng Xu. Arbitrary style transfer via multi-adaptation network. In *Proceedings of the 28th ACM international conference on multimedia*, 2020. 2

[13] Yu Deng, Jiaolong Yang, and Xin Tong. Deformed implicit field: Modeling 3d shapes with learned dense correspondence. In *CVPR*, 2021. 3

[14] Zhiwen Fan, Yifan Jiang, Peihao Wang, Xinyu Gong, Dejia Xu, and Zhangyang Wang. Unified implicit neural stylization. In *ECCV*, 2022. 2

[15] Anna Frühstück, Nikolaos Sarafianos, Yuanlu Xu, Peter Wonka, and Tony Tung. VIVE3D: Viewpoint-independent video editing using 3D-aware GANs. In *CVPR*, 2023. 2

[16] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *ICCV*, 2021. 3

[17] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR*, 2016. 1, 2

[18] Xinyu Gong, Haozhi Huang, Lin Ma, Fumin Shen, Wei Liu, and Tong Zhang. Neural stereoscopic image style transfer. In *ECCV*, 2018. 2

[19] Shuyang Gu, Congliang Chen, Jing Liao, and Lu Yuan. Arbitrary style transfer with deep feature reshuffle. In *CVPR*, 2018. 2

[20] Ayaan Haque, Matthew Tancik, Alexei Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *ICCV*, 2023. 3

[21] Eric Heitz, Kenneth Vanhoey, Thomas Chambon, and Laurent Belcour. A sliced wasserstein loss for neural texture synthesis. In *CVPR*, 2021. 2

[22] Douglas R. Hofstadter. Metamagical themas: Variations on a theme as the essence of imagination. *Scientific American*, 1983. 2

[23] Lukas Höllein, Justin Johnson, and Matthias Nießner. Stylemesh: Style transfer for indoor 3d scene reconstructions. In *CVPR*, 2022. 1, 2

[24] 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*, 2017. 1, 2

[25] Hsin-Ping Huang, Hung-Yu Tseng, Saurabh Saini, Maneesh Singh, and Ming-Hsuan Yang. Learning to stylize novel views. In *ICCV*, 2021. 1, 2

[26] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017. 2

[27] Zixuan Huang, Jinghuai Zhang, and Jing Liao. Style mixer: Semantic-aware multi-style transfer network. In *Computer Graphics Forum*. Wiley Online Library, 2019. 2

[28] Yongcheng Jing, Yang Liu, Yezhou Yang, Zunlei Feng, Yizhou Yu, Dacheng Tao, and Mingli Song. Stroke controllable fast style transfer with adaptive receptive fields. In *ECCV*, 2018. 5

[29] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 2

[30] Yucheol Jung, Wonjong Jang, Soongjin Kim, Jiaolong Yang, Xin Tong, and Seungyong Lee. Deep deformable 3d caricatures with learned shape control. In *ACM SIGGRAPH 2022 Conference Proceedings*, 2022. 3

[31] Hiromichi Kamata, Yuiko Sakuma, Akio Hayakawa, Masato Ishii, and Takuya Narihira. Instruct 3d-to-3d: Text instruction guided 3d-to-3d conversion. *arXiv preprint arXiv:2303.15780*, 2023. 3

[32] Zhanghan Ke, Chunyi Sun, Lei Zhu, Ke Xu, and Rynson WH Lau. Harmonizer: Learning to perform white-box image and video harmonization. In *ECCV*, 2022. 2

[33] Sunnie SY Kim, Nicholas Kolkin, Jason Salavon, and Gregory Shakhnarovich. Deformable style transfer. In *ECCV*, 2020. 2

[34] Nicholas Kolkin, Jason Salavon, and Gregory Shakhnarovich. Style transfer by relaxed optimal transport and self-similarity. In *CVPR*, 2019. 1, 2

[35] Nicholas Kolkin, Michal Kucera, Sylvain Paris, Daniel Sykora, Eli Shechtman, and Greg Shakhnarovich. Neural neighbor style transfer. *arXiv preprint arXiv:2203.13215*, 2022. 3, 6

[36] Dmytro Kotovenko, Artsiom Sanakoyeu, Sabine Lang, and Bjorn Ommer. Content and style disentanglement for artistic style transfer. In *ICCV*, 2019. 2

[37] Dmytro Kotovenko, Artsiom Sanakoyeu, Pingchuan Ma, Sabine Lang, and Bjorn Ommer. A content transformation block for image style transfer. In *CVPR*, 2019.

[38] Dmytro Kotovenko, Matthias Wright, Arthur Heimbrecht, and Bjorn Ommer. Rethinking style transfer: From pixels to parameterized brushstrokes. In *CVPR*, 2021. 2

[39] Dishani Lahiri, Neeraj Panse, and Moneish Kumar. S2rf: Semantically stylized radiance fields. In *ICCVW*, 2023. 8

[40] Chuan Li and Michael Wand. Combining markov random fields and convolutional neural networks for image synthesis. In *CVPR*, 2016. 1, 2

[41] Xueting Li, Sifei Liu, Jan Kautz, and Ming-Hsuan Yang. Learning linear transformations for fast arbitrary style transfer. *arXiv preprint arXiv:1808.04537*, 2018. 1, 2

[42] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. In *NeurIPS*, 2017. 2

[43] Jing Liao, Yuan Yao, Lu Yuan, Gang Hua, and Sing Bing Kang. Visual attribute transfer through deep image analogy. *ACM TOG*, 36(4), 2017. 2

[44] Hsueh-Ti Derek Liu, Michael Tao, and Alec Jacobson. Paparazzi: Surface editing by way of multi-view image processing. *ACM Transactions on Graphics*, 2018. 1, 2

[45] Kunhao Liu, Fangneng Zhan, Yiwen Chen, Jiahui Zhang, Yingchen Yu, Abdulmotaleb El Saddik, Shijian Lu, and Eric P Xing. Stylerf: Zero-shot 3d style transfer of neural radiance fields. In *CVPR*, 2023. 1, 2, 3

[46] Songhua Liu, Tianwei Lin, Dongliang He, Fu Li, Meiling Wang, Xin Li, Zhengxing Sun, Qian Li, and Errui Ding. Adaattn: Revisit attention mechanism in arbitrary neural style transfer. In *ICCV*, 2021. 2

[47] Xiao-Chang Liu, Xuan-Yi Li, Ming-Ming Cheng, and Peter Hall. Geometric style transfer. *arXiv preprint arXiv:2007.05471*, 2020. 2

[48] Xiao-Chang Liu, Yong-Liang Yang, and Peter Hall. Learning to warp for style transfer. In *CVPR*, 2021. 2

[49] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. Deep photo style transfer. In *CVPR*, 2017. 1, 2

[50] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. Deep painterly harmonization. In *Computer graphics forum*. Wiley Online Library, 2018. 2

[51] Roey Mechrez, Itamar Talmi, and Lihi Zelnik-Manor. The contextual loss for image transformation with non-aligned data. In *ECCV*, 2018. 2

[52] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38, 2019. 5, 6

[53] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1, 2

[54] Fangzhou Mu, Jian Wang, Yicheng Wu, and Yin Li. 3d photo stylization: Learning to generate stylized novel views from a single image. In *CVPR*, 2022. 1, 2

[55] Thu Nguyen-Phuoc, Feng Liu, and Lei Xiao. Snerf: stylized neural implicit representations for 3d scenes. *ACM TOG*, 41, 2022. 1, 2, 5, 6, 7

[56] Hong-Wing Pang, Binh-Son Hua, and Sai-Kit Yeung. Locally stylized neural radiance fields. In *ICCV*, 2023. 1, 2

[57] Yingxue Pang, Jianxin Lin, Tao Qin, and Zhibo Chen. Image-to-image translation: Methods and applications. *IEEE Transactions on Multimedia*, 24, 2021. 5

[58] Dae Young Park and Kwang Hee Lee. Arbitrary style transfer with style-attentional networks. In *CVPR*, 2019. 2

[59] 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*, 2021. 3

[60] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, 2021. 3

[61] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, 2021. 2, 3

[62] Eric Risser, Pierre Wilmot, and Connelly Barnes. Stable and controllable neural texture synthesis and style transfer using histogram losses. *arXiv preprint arXiv:1701.08893*, 2017. 2

[63] Manuel Ruder, Alexey Dosovitskiy, and Thomas Brox. Artistic style transfer for videos and spherical images. *IJCV*, 126(11), 2018. 2

[64] Artsiom Sanakoyeu, Dmytro Kotovenko, Sabine Lang, and Bjorn Ommer. A style-aware content loss for real-time hd style transfer. In *ECCV*, 2018. 2

[65] Jonathan Shade, Steven Gortler, Li-wei He, and Richard Szeliski. Layered depth images. In *Proceedings of the 25th Annual Conference on Computer Graphics and Interactive Techniques*, 1998. 5

[66] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. Singan: Learning a generative model from a single natural image. In *ICCV*, 2019. 5, 7

[67] Lu Sheng, Ziyi Lin, Jing Shao, and Xiaogang Wang. Avatarnet: Multi-scale zero-shot style transfer by feature decoration. In *CVPR*, 2018. 2

[68] Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Bulò, Norman Müller, Matthias Nießner, Angela Dai, and Peter Kontschieder. Panoptic lifting for 3d scene understanding with neural fields. In *CVPR*, 2023. 2, 8

[69] K Simonyan and A Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 2, 3, 5

[70] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Nonrigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *ICCV*, 2021. 3

[71] Yi-Hsuan Tsai, Xiaohui Shen, Zhe Lin, Kalyan Sunkavalli, Xin Lu, and Ming-Hsuan Yang. Deep image harmonization. In *CVPR*, 2017. 2

[72] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *CVPR*, 2022. 3

[73] Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Nerf-art: Text-driven neural radiance fields stylization. *IEEE Transactions on Visualization and Computer Graphics*, 2023. 3

[74] Weiyue Wang, Duygu Ceylan, Radomir Mech, and Ulrich Neumann. 3dn: 3d deformation network. In *CVPR*, 2019. 3

[75] Wenjing Wang, Jizheng Xu, Li Zhang, Yue Wang, and Jiaying Liu. Consistent video style transfer via compound regularization. *AAAI*, 2020. 2

[76] Wenjing Wang, Shuai Yang, Jizheng Xu, and Jiaying Liu. Consistent video style transfer via relaxation and regularization. *IEEE Transactions on Image Processing*, 29, 2020. 1, 2

[77] Xiaolei Wu, Zhihao Hu, Lu Sheng, and Dong Xu. Styleformer: Real-time arbitrary style transfer via parametric style composition. In *ICCV*, 2021. 2

[78] Zijie Wu, Zhen Zhu, Junping Du, and Xiang Bai. Ccpl: contrastive coherence preserving loss for versatile style transfer. In *ECCV*. Springer, 2022. 5

[79] Xide Xia, Meng Zhang, Tianfan Xue, Zheng Sun, Hui Fang, Brian Kulis, and Jiawen Chen. Joint bilateral learning for real-time universal photorealistic style transfer. In *ECCV*, 2020. 2

[80] Xide Xia, Tianfan Xue, Wei-sheng Lai, Zheng Sun, Abby Chang, Brian Kulis, and Jiawen Chen. Real-time localized photorealistic video style transfer. 2021. 1, 2

[81] Shuai Yang, Zhangyang Wang, Zhaowen Wang, Ning Xu, Jiaying Liu, and Zongming Guo. Controllable artistic text style transfer via shape-matching gan. In *ICCV*, 2019. 2

[82] Jordan Yaniv, Yael Newman, and Ariel Shamir. The face of art: Landmark detection and geometric style in portraits. *ACM TOG*, 38(4), 2019. 2

[83] Kai Zhang, Nick Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snavely. Arf: Artistic radiance fields. In *ECCV*, 2022. 1, 2, 3, 5, 6, 7

[84] Lingzhi Zhang, Tarmily Wen, and Jianbo Shi. Deep image blending. In *WACV*, 2020. 2

[85] Pan Zhang, Bo Zhang, Dong Chen, Lu Yuan, and Fang Wen. Cross-domain correspondence learning for exemplar-based image translation. In *CVPR*, 2020. 2

[86] Yuechen Zhang, Zexin He, Jinbo Xing, Xufeng Yao, and Jiaya Jia. Ref-NPR: Reference-based non-photorealistic radiance fields for controllable scene stylization. In *CVPR*, 2023. 1, 2, 3, 5, 6, 7

[87] Zicheng Zhang, Yinglu Liu, Congying Han, Yingwei Pan, Tiande Guo, and Ting Yao. Transforming radiance field with lipschitz network for photorealistic 3d scene stylization. In *CVPR*, 2023. 2, 3