

\mathcal{G} -Style: Stylized Gaussian Splatting

Áron Samuel Kovács , Pedro Hermosilla, and Renata G. Raidou 

TU Wien, Austria

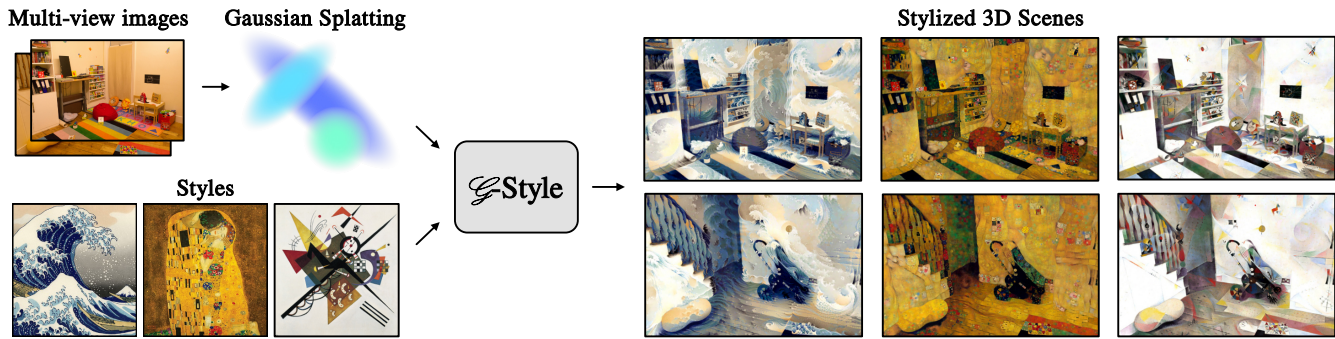


Figure 1: \mathcal{G} -Style: Our method takes a 3D scene, represented using Gaussian Splatting, and a style image exemplar as input, and generates a stylized version of the scene that closely matches the visual style of the exemplar. By modifying the geometry of the scene and designing losses that capture style patterns at different scales, we achieve high-quality stylized scenes efficiently, with results generated in just a few minutes.

Abstract

We introduce \mathcal{G} -Style, a novel algorithm designed to transfer the style of an image onto a 3D scene represented using Gaussian Splatting. Gaussian Splatting is a powerful 3D representation for novel view synthesis, as—compared to other approaches based on Neural Radiance Fields—it provides fast scene renderings and user control over the scene. Recent pre-prints have demonstrated that the style of Gaussian Splatting scenes can be modified using an image exemplar. However, since the scene geometry remains fixed during the stylization process, current solutions fall short of producing satisfactory results. Our algorithm aims to address these limitations by following a three-step process: In a pre-processing step, we remove undesirable Gaussians with large projection areas or highly elongated shapes. Subsequently, we combine several losses carefully designed to preserve different scales of the style in the image, while maintaining as much as possible the integrity of the original scene content. During the stylization process and following the original design of Gaussian Splatting, we split Gaussians where additional detail is necessary within our scene by tracking the gradient of the stylized color. Our experiments demonstrate that \mathcal{G} -Style generates high-quality stylizations within just a few minutes, outperforming existing methods both qualitatively and quantitatively.

CCS Concepts

• Computing methodologies → Computer graphics; Artificial intelligence; Neural networks;

1. Introduction

While humans excel at creating paintings with specific contents and styles, this task has proven challenging for computers to replicate. With the advent of neural networks—and in particular, Convolutional Neural Networks—algorithms have been developed to transfer the style of one image onto another [GEB15a]. In this process, a *content image* refers to the original image whose subject matter we aim to retain, while a *style image* is the one whose artistic style we want to apply to the content image. These algorithms enabled the modification of the style of an

image while preserving its content by matching the statistical properties of the embeddings of both content and style images, as obtained from a deep neural network.

With the appearance of novel view synthesis methods based on neural networks (NeRFs) [MST*20], researchers turned their attention to applying style transfer techniques to entire 3D scenes. Style transfer for 3D scenes aims to generate novel views of a scene from a finite number of images of the same scene with a particular style specified by a style image exemplar. To succeed in this task, the style transfer

method should ensure multi-view consistency between views to provide a smooth navigation experience. Despite the high-quality results provided by 2D style transfer methods, the same algorithms were not able to provide *consistent styles* across 3D scene views [HTS*21, HHY*22, MWL22, NPLX22]. Therefore, they have been deemed unfit for 3D style transfer. The limitations of these works have been addressed by several methods [HHY*22, NPLX22, ZKB*22a]—either by modifying the colors of the pre-trained NeRF representations or by techniques such as color transfer [ZKB*22a], provided that the original methods were able to support view-dependent effects. Still, NeRF-based approaches require *large training* and *rendering times*.

Recently, Gaussian Splatting (GS) has established itself as an active area of research [CW24, FXZ*24, WYZ*24] by demonstrating that it can excel in view synthesis quality while requiring significantly less time for optimization and rendering. Two recent pre-prints [LZX*24, SGC*24] have suggested using GS as the main scene representation and modifying the color of these to obtain stylized renderings of the scene. These methods do not alter the position and shape of the Gaussian representing the scene, which results in a *low resolution* in some areas of the scene and, hence, *low-resolution stylized colors*. In addition to the aforementioned problems, all 3D style transfer methods only focus on transferring *high-frequency patterns* from the style image, such as brush strokes or color statistics. However, the notion of the style of an image is more nuanced and can significantly differ from image to image. In some cases, like in Figure 1, the style is mostly defined by large patterns that existing methods might miss—thus, reducing their effectiveness in transferring the intended style.

To overcome all these limitations, we introduce \mathcal{G} -Style, a novel method to transfer the style of an image onto a 3D scene represented using Gaussian Splatting. Our approach requires only a few minutes to optimize and delivers a high-quality stylized GS representation of the scene. Our method comprises several steps: First, we pre-process the scene represented with a set of Gaussians to ensure a uniform coverage. Then, our method starts the stylization process by updating the color associated with each Gaussian. During the stylization process, we enhance the resolution by splitting Gaussians with high color gradients, adding finer details where necessary. Due to the sparse nature of the scene and by only representing diffuse surfaces, our method provides view consistency by construction. Moreover, our dual loss function enables the algorithm to capture both high-frequency and low-frequency patterns in the style image, therefore generating high-quality 3D scene renderings that reproduce a large variety of artistic styles. Our extensive evaluation demonstrates that our approach outperforms existing methods in style transfer quality and rendering time.

2. Related Work

2D Style Transfer. Gatys et al. [GEB15a] first introduced a method for neural style transfer, where the artistic style of a provided image is used to visually reconstruct the content of another image. The method is based on iterative optimization to match the output and second-order statistics, expressed as the Gram matrices of hidden layers of a pre-trained network. If preserving the content is not necessary, the method can synthesize a texture resembling the provided style image, initializing the process with a noise image [GEB15b] and disregarding the part of the original loss function that is responsible for content preservation.

This approach has been since refined to enable transferring features

on multiple scales [ZGW*22], by utilizing Generative Adversarial Networks (GANs) in a course-to-fine fashion [JBV17] or diffusion models [ZHT*23, WZX23, CHH24]. Also, different ways of capturing and matching the style statistics have been proposed [LYY*17, GCLY18, KSS19]. All aforementioned methods aim to improve style details by introducing pen or brush strokes, which were too blurry in the original work of Gatys et al., or by better preserving semantic consistency.

Instead of relying on matching extracted style statistics, searching for nearest neighbors in the feature space and minimizing distances between them is another option for transferring style [KSS19, CS16, LYY*17, LW16, ZKB*22a]. This strategy can lead to significant improvements in transferring high-frequency features, by avoiding the averaging of features that match possibly multi-modal style statistics.

Existing methods can be further subdivided based on whether they iteratively modify an image, just like in the original work of Gatys et al., or explicitly minimize an objective function with a single feed-forward pass [CS16, HB17, AHS*21]. While the feed-forward-based approaches are orders of magnitude faster, their results are generally of lower quality than the slower optimization-based techniques.

3D Style Transfer. 3D style transfer refers to modifying the appearance of a 3D object or a scene so that when viewed from different angles, it matches the style of a given exemplar. To reuse a 2D style transfer method, one needs to extract and employ a 2D image of a given scene. This can be either obtained by utilizing a differentiable renderer [MPSO18, ZKB*22a], by slicing the 3D volume [HMR20, GRGH19, ZGW*22, CW10, KFCO*07], or by directly working on the surface manifold of textured meshes [KHR24].

Existing approaches usually rely on a single way of representing objects and whole scenes, which is also the main point of influence for their performance. The approach of Cao et al. [CWNN20] uses point clouds, yielding holes due to measurement errors during scanning and being unsuitable for representing surfaces and real-world scenes. Conversely, Mordvintsev et al. [MPSO18] and Hoellein et al. [HJN22] use textured meshes, which can be reconstructed from the aforementioned point clouds. Due to the discrete and restrictive nature of mesh representations, these approaches can suffer from different types of artifacts.

Lately, Neural Radiance Fields (NeRFs) [MST*20] and Gaussian Splatting [KKLD23] have become very prominent in reconstructing objects and scenes. Unlike point clouds and meshes, NeRFs and Gaussians are soft volumetric representations. Naturally, these representations can be utilized for scene reconstruction or style transfer. Like in the 2D case, these style transfer methods can be subdivided into two categories: zero-shot and iterative methods. Zero-shot methods [LZL*23, LZC*23a, XCX*24] focus on matching colors, relighting, or transferring details on a small scale to achieve multi-view consistency. This limitation arises from their inability to quickly and consistently embed large features into scenes. Thus, these methods struggle to synthesize large patterns that span significant portions of a given scene.

On the other hand, iterative methods leverage rendering scenes from multiple viewpoints. Hence, they can construct large-scale patterns that remain consistent across different views, while any inconsistencies are corrected during the training process. Zhang et al. [ZKB*22a] propose a new loss based on nearest neighbor feature matching (NNFM) which better preserves details and also optimizes backpropagation by deferring the changes to the network after first computing the losses for full-resolution

images. Zhang et al. [ZFLS24] further improve this work by enabling the stylization of specific objects and introducing a semantic-aware version of NNFM. These methods suffer from limited user control and require significant time to optimize—limitations that can be addressed by a novel view synthesis alternative: Gaussian Splatting [KKLD23].

In this work, we focus on style transfer in scenes represented with the help of 3D Gaussian Splatting, allowing us to obtain a stylized scene representation in a few minutes and supporting real-time inspection of the result. Two concurrent methods address the same problem as we do from different angles. The first, StyleGaussian, embeds 2D VGG scene features into reconstructed 3D Gaussians, transforms them according to a style, and decodes them onto a stylized image [LZX*24]. The second, GSS, employs pre-trained Gaussians conditioned on a style image to obtain stylized views of complex 3D scenes with spatial consistency [SGC*24]. However, these methods do not change the geometry of the scene during optimization resulting in low-resolution stylized colors. As we demonstrate later, our method achieves higher-quality scene representations that preserve large patterns from the style image.

3. Background: NeRFs and Gaussian Splatting

In this section, we briefly describe the two most commonly used 3D representations for novel view synthesis, NeRFs [MST*20] and Gaussian Splatting [KKLD23].

NeRFs. Neural radiance fields have revolutionized the field of novel view synthesis by introducing a new scene representation, and an optimization algorithm to train this representation only from images. The outgoing radiance at any point x in the scene for any view direction v is modeled by a parametric model $\phi_{\theta}(x, v)$ with parameters θ . This model is usually a Multi-layer Perceptron (MLP), which is chosen due to the universality provided by this type of model. To train this model, the scene is rendered from multiple views using the volume rendering algorithm [Max95]. By comparing the generated image \mathcal{I}_{∇} to a real picture of the scene \mathcal{I}_{gt} , gradients can be computed for the parameters θ through the rendering operation and the scene representation can be updated to match the ground truth images. Despite the high-quality images generated by NeRFs, the control provided to the user is limited, and they demand extensive rendering times due to the multiple evaluations required to compute the value of a pixel.

Gaussian Splatting. Gaussian Splatting [KKLD23] has emerged as a viable alternative to address the main limitations of NeRFs: limited scene control and low rendering speed. To reconstruct a real-life scene from ground truth images \mathcal{I}_{gt} , Kerbl et al. use the Structure from Motion algorithm [Ull79] to obtain camera poses for each image and a sparse set of initial 3D points. The 3D points are then transformed into Gaussian functions, each representing a point in the scene. In this way, Gaussian Splatting represents the function $\phi_{\theta}(x, v)$ as a set of 3D Gaussians.

Each Gaussian is defined by its mean μ , covariance matrix Σ , opacity δ , and color c . By representing the color with spherical harmonics, view-dependent effects can also be captured. Since the covariance matrix Σ has to be positive semi-definite, which is difficult to enforce during optimization, the covariance matrix is obtained by $RSS^T R^T$, where S is a scaling matrix and R a rotation matrix obtained from a quaternion. The covariance matrices of the Gaussians are initialized accounting for their neighbors to conservatively cover the surfaces and prevent holes in the reconstruction.

The optimization process used in Gaussian Splatting is similar to the one used in NeRFs, where volume rendering is used to generate images by querying density and color along the rays emanating from the camera. However, due to the sparse nature of the representation achieved with Gaussian Splatting, the rendering process can be efficiently computed by projecting the Gaussians into the image and combining them using alpha blending. Furthermore, since the initial set of Gaussians may not be sufficient to capture all the necessary geometric and color details, the Gaussians are periodically split or cloned based on their accumulated μ gradient. However, this splitting is designed to represent the original scene and might not be sufficient for its stylized version. As we discuss in the upcoming sections, we propose a fast, high-quality, and consistent approach for stylizing 3D scenes by modifying the style and geometry of scenes represented by Gaussian Splatting with the style of an additional image or texture.

4. Methodology of \mathcal{G} -Style: Gaussian Splatting with Style

In this section, we describe our proposed algorithm: \mathcal{G} -Style. We first provide an overview of our approach (also illustrated in Figure 2), followed by a detailed explanation of its substeps.

4.1. Overview

Our algorithm takes as input a scene represented with a set of Gaussians \mathcal{G} and a set of ground truth images \mathcal{I}_{gt} . Subsequently, it modifies \mathcal{G} based on the style provided from a style exemplar \mathcal{I}_s . First, we *pre-process* \mathcal{G} to remove long narrow Gaussians and Gaussians covering large areas, making the initial set of Gaussians uniform. Once the initial representation has been pre-processed, we create an additional color c_s associated with each Gaussian which is initialized with the original color c_{gt} . These new colors c_s are modified during a *stylization* process that uses a composition of several losses to preserve different properties of the style of \mathcal{I}_s . Since the geometry provided by the initial pre-processing step can be limited to represent detailed style features, \mathcal{G} undergoes a *geometric fine-tuning* step. In this step, the Gaussians are split based on the gradient of c_s and fine-tuned to match the original scene images \mathcal{I}_{gt} by modifying μ , Σ , δ , and c_{gt} . The stylization and geometric fine-tuning steps are repeated until convergence.

4.2. Pre-processing Step

Although Gaussian Splatting offers high-quality scene representation, the original approach has limitations. It often generates large flat Gaussians to depict uniform, flat surfaces (e.g., walls), and narrow elongated Gaussians to capture high-frequency details. The latter might also come as a byproduct of the optimization process. When incorporating an additional style into the scene, large flat areas might need additional Gaussians to incorporate extra details, while already highly detailed areas may not need more Gaussians to stylize them properly. Therefore, in the initial step of our approach, the scene undergoes a Gaussian normalization process where the resulting representation \mathcal{G} is composed of Gaussians of similar size and shape.

Flat Gaussian Split. To detect under-sampled areas, we compute the approximated maximum projected area of each Gaussian, A_i . We compute A_i by multiplying the two highest components of the Gaussian’s scaling matrix S_i . Then, we mark for splitting all Gaussian above a threshold $t_f = \mu_A + \gamma\sigma_A$, where μ_A is the mean of all A in

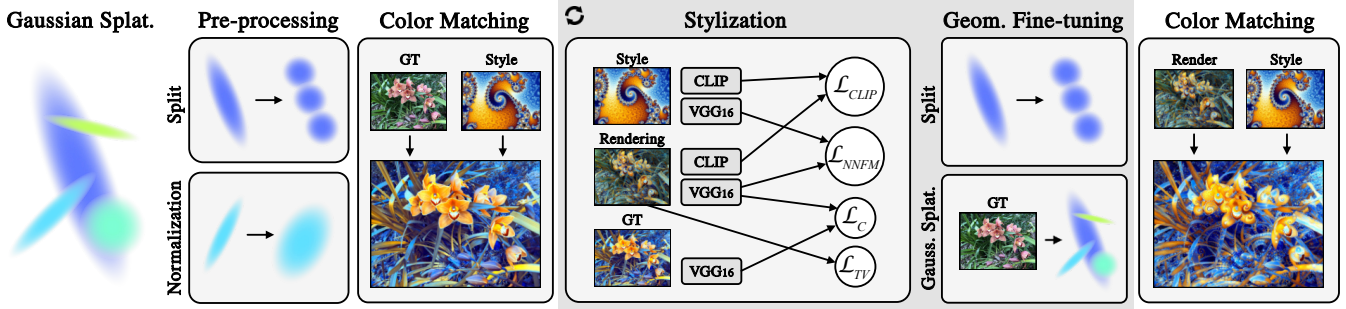


Figure 2: Overview of our method: We take a 3D scene represented using Gaussian Splatting and pre-process it to subdivide large Gaussians and normalize elongated ones. Initially, we perform a color matching between the ground truth images and the style image. Subsequently, we start the iterative stylization process. First, we optimize the colors of the Gaussians using multiple losses to capture style patterns at different scales while preserving the content of the scene. Then, we fine-tune the geometry of the scene and add details for Gaussians with a large gradient of the stylized color. We repeat the stylization and geometry fine-tuning steps until convergence. At the end, we perform an additional color matching step between the renderings of the resulting scene and the style image.

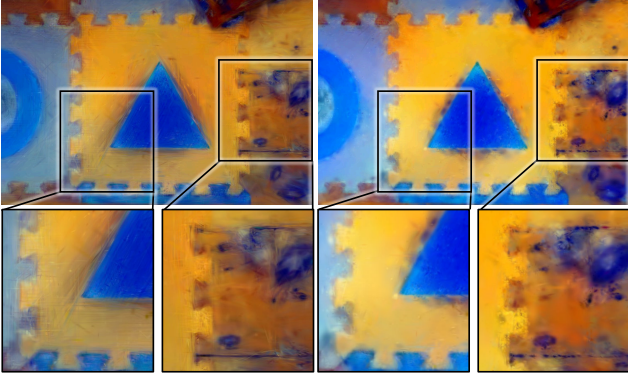


Figure 3: Gaussian Normalization: We normalize (right) the size of narrow Gaussians to avoid multiple overlying Gaussians (left).

the scene, σ_A is the standard deviation of A , and γ is a user-defined parameter controlling the number of Gaussian mark for splitting. As splitting Gaussians modifies the geometry of the scene, we are obliged to optimize them. For this, we employ the same optimization algorithm as in the original work of Kerbl et al. [KKLD23].

Narrow Gaussian Normalization. To identify Gaussians with a narrow shape, we compute their elongation factor E_i by dividing the highest component of each Gaussian’s scaling matrix S_i by its second highest component. If E_i is above a certain user-defined threshold t_e , we mark it for normalization, which sets the largest component to the average of the largest and second-largest components. We repeatedly perform this operation during the optimization process after the flat Gaussian split. This optimization process retrains the scene to match the appearance of the ground truth images \mathcal{I}_{gt} . As such, it corrects the deformations caused by this narrowing step while keeping the Gaussians more rounded. In Figure 3, we illustrate the effect of normalizing Gaussians (right) as opposed to not normalizing them (left).

Diffuse Color Transform. In 3D style transfer, where there is



Figure 4: Pre-processing: The effect of our pretraining step. Left: before pretraining, right: after pretraining. The scale of Gaussians is set to 0.25, otherwise both images would look identical.

no concrete ground truth image associated with a view, enforcing multi-view consistency is key for seamless navigation through the scene. In the original Gaussian Splatting algorithm, spherical harmonics enable modeling view-dependent effects such as reflections. However, they also generate undesirable artifacts on the stylized version of the scene, where each view direction could result in a completely different stylization. To avoid this, we only use the zeroth term of the spherical harmonics representation, limiting the model to represent diffuse objects only, therefore, enforcing view consistency.

Parameterizations. We perform five rounds of the splitting–normalization–optimization process, which in our experiments offers a good balance between the pre-computation time and the even distribution of Gaussian sizes. Initially, we set γ to 1.1, but in each subsequent round, we multiply it by 1.125 to deal with very large Gaussians that may still be remaining. The initial value for γ is quite low and, thus, we split at most 5% of the largest Gaussian in each round. Additionally, we set the threshold for elongation t_e to 1.5, which allows a degree of elongation but normalizes the shape of those Gaussians that heavily affect the final appearance. After pre-processing, the resulting set of Gaussians still matches the visual characteristics of a scene, but the distribution of their perceived sizes is uniform, as can be seen in Figure 4. Note that incorporating more Gaussians into the scene results in additional details—and additional memory demands. However, since

we limit our representation to diffuse objects, the storage required to store the color of each Gaussian is reduced from 192 bytes to 12 bytes for the three components of the diffuse color. Our implementation, thus, allows us to increase the number of Gaussians to represent the scene while at the same time reducing memory consumption.

4.3. Stylization Step

Once we have pre-processed the scene, we start the stylization process. During stylization, we render the scene from multiple views using c_s and compute several losses over the images carefully designed to preserve the style of \mathcal{I}_s at different scales. We hereby describe in detail the different losses used in our algorithm.

Low-frequency Style. The style of an image is determined by color and high-frequency details and also by large-scale patterns and features. To preserve these low-frequency features of the style image, we employ a CLIP-based [RKH*21] loss. We encode our rendered images \mathcal{I}_r^k and the style image \mathcal{I}_s into a feature vector using the image encoder of the CLIP model, \mathcal{C} . Then, our loss is defined as the similarity between these feature vectors where similarity is measured as the l_2 :

$$\mathcal{L}_{CLIP} = \frac{1}{K} \sum_{k=1}^K \|\mathcal{C}(\mathcal{I}_r^k) - \mathcal{C}(\mathcal{I}_s)\|_2^2 \quad (1)$$

where K is the number of rendered images in the batch.

High-frequency Style. Our CLIP-based loss can capture large-scale patterns in the style image. Yet, fine details, such as brush strokes in a painting, might not be well represented with this loss. Therefore, following Zhang et al. [ZKB*22a], we use a nearest neighbor feature matching (NNFM) loss utilizing a pre-trained VGG-16 network [SZ14] to capture the high-frequency style patterns. The architecture of a feature extractor can play a significant role in the quality of the generated results. We selected VGG-16 for its proven effectiveness in previous style transfer solutions, facilitating easier comparisons with those methods. The NNFM loss is a replacement for the widely used Gram matrix-based loss of Gatys et al. [GEB15a]. Instead of matching statistics of feature maps \mathcal{F}_r and \mathcal{F}_s , this loss searches for nearest neighbors in the feature space. Let \mathcal{F}_r and \mathcal{F}_s be the VGG-16 features maps of \mathcal{I}_r and \mathcal{I}_s respectively, and let $\mathcal{F}(i, j)$ be the feature vector at pixel location (i, j) . The NNFM loss is then defined as:

$$\mathcal{L}_{NNFM}(\mathcal{F}_r, \mathcal{F}_s) = \frac{1}{N} \sum_{i,j} \min_{i',j'} D(\mathcal{F}_r(i, j), \mathcal{F}_s(i', j')) \quad (2)$$

where N is the number of pixels in \mathcal{F}_r , and D is the cosine distance between two vectors. As in Zhang et al. [ZKB*22a], we use feature maps from the third block of VGG-16.

Regularization. To avoid the deterioration of features that are necessary for scene understanding, like in the work of Zhang et al. [ZKB*22a], we utilize a content loss \mathcal{L}_C , which is the l_2 loss between the features of rendered images \mathcal{F}_r and ground truth images \mathcal{F}_{gr} . This loss also uses the features from the third block of VGG-16. Additionally, we incorporate a total variation term in our loss \mathcal{L}_{TV} to prevent noise in the resulting renderings.

Complete Style Loss. Our final loss is a weighted sum of all the losses, given by the equation:

$$\mathcal{L} = \lambda_{CLIP} \mathcal{L}_{CLIP} + \lambda_{NNFM} \mathcal{L}_{NNFM} + \lambda_C \mathcal{L}_C + \lambda_{TV} \mathcal{L}_{TV} \quad (3)$$

We perform the stylization process for 15 epochs. Note that forward-facing scenes do not have as strict multi-view consistency requirements as 360° scenes due to the limited viewing angles of the ground truth images. As such, they converge more easily compared to 360° scenes. To account for the differences between these two types of scenes, we use different parameterizations. For forward-facing scenes, we use an exponentially decaying learning rate from 1e-1 to 1e-2, and for 360° scenes the learning rate decays from 1e-2 to 5e-3. Also, we set λ_{CLIP} to 10, λ_{NNFM} to 100, λ_C to 0.05, and λ_{TV} to 1e-4. For 360° scenes, we set λ_{NNFM} to 10.

4.4. Geometric Fine-tuning Step

In the pre-processing step, we place more Gaussians in the undersampled areas, making the Gaussians more uniform in size. This step is conservative enough to not overly increase the number of Gaussians and, thus, not oversample a given scene. With the new Gaussians, it is possible to synthesize features that otherwise could not be represented. However, the size of the Gaussians still limits the synthesis of very fine features in our stylized scene.

To overcome this, during the stylization process, we periodically split Gaussians based on their c_s gradient. Similarly to the work of Kerbl et al. [KKLD23], we keep an accumulation buffer \mathcal{B} to store the norm of the gradient c_s . After each iteration and for each Gaussian, we add the norm of the c_s gradient to \mathcal{B} and periodically split a user-defined percentage of Gaussians with the highest value. Since splitting Gaussians modifies the geometry of the scene, after each splitting, we optimize their μ , Σ , δ , and c_{gr} in the same optimization process as Kerbl et al. [KKLD23].

4.5. Style Color Matching

To ensure a similar color distribution between the resulting stylized 3D scene and the original style image, at the beginning of our algorithm, we match the mean and covariance matrix of colors from our ground truth images \mathcal{I}_{gr} and the style image \mathcal{I}_s . Let C be a matrix containing pixel colors of the ground truth images of the scene \mathcal{I}_{gr} and S be a matrix of pixels from \mathcal{I}_s , where each row is one pixel and columns are used for RGB components. We analytically solve for a linear transformation A , such that $E(AC) = E(S)$ and $Cov(AC) = Cov(S)$, and finally modify our initial Gaussian Splatting scene representation so that the rendered color images match with the color-corrected ground truth images. Even though this color transfer step assumes unimodal distributions of colors, in our experiments, we empirically identified that it is sufficient for all tested style images. Since optimizing with a loss that considers the activations of hidden layers does not guarantee that the resulting colors are accurate to the style, we apply the same correction at the end of our algorithm to \mathcal{I}_r .

5. Results

In this section, we provide an analysis of the results produced by our method. Additionally, we offer a comparison to other leading state-of-the-art approaches. All of our results are generated using a modified 3D Gaussian Splatting codebase [Ker23], which is publicly available in our GitHub Repository (<https://github.com/AronKovacs/g-style>).

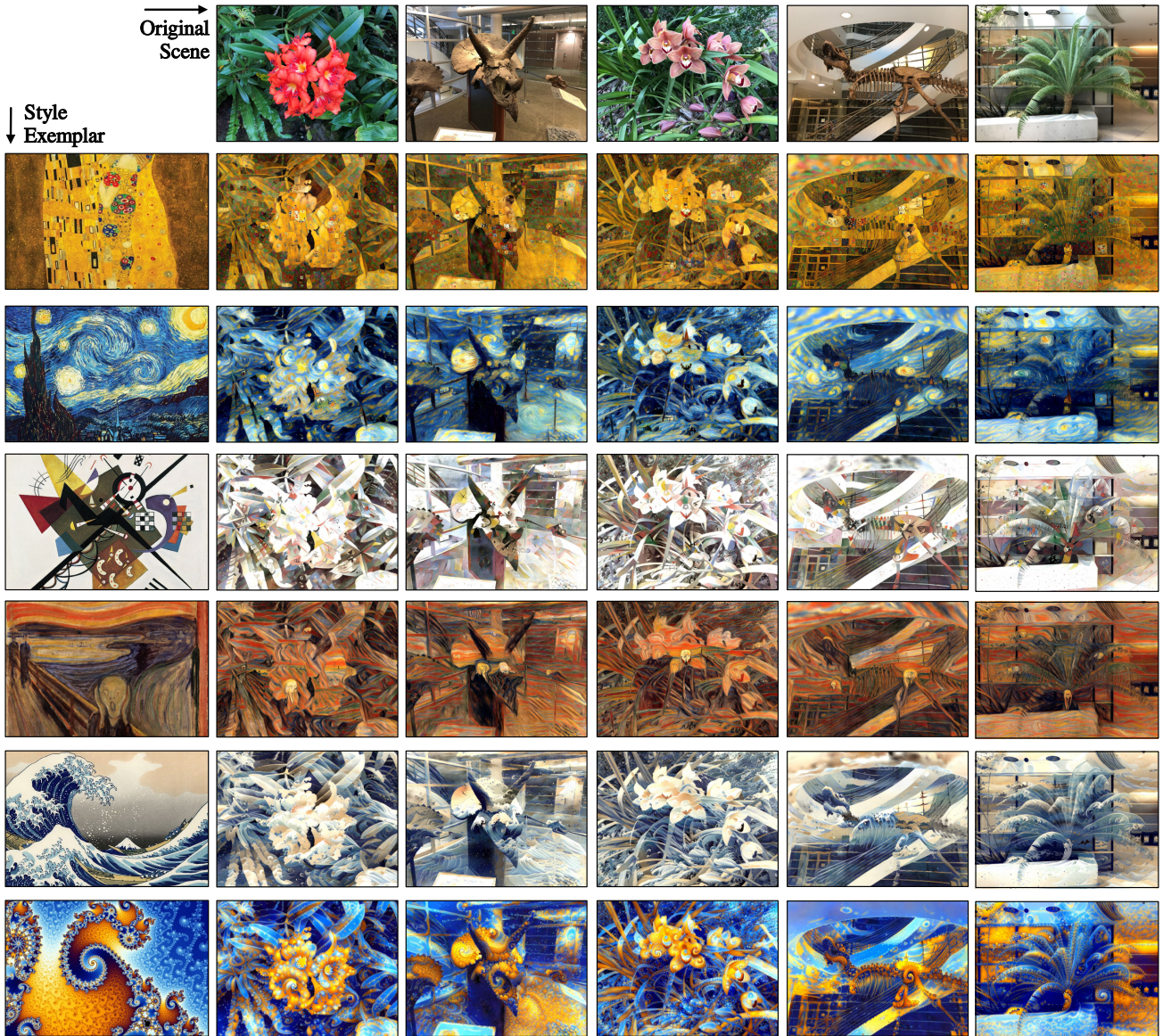


Figure 5: Results generated with our approach, \mathcal{G} -Style, for five forward-facing scenes (columns) given six style exemplars (rows).

5.1. Datasets

To evaluate our approach, we prepared a series of forward-facing scenes: *Flower*, *Horns*, *Orchid*, *T-Rex*, and *Fern* (employed in the work of Mildenhall et al. [MSOC*19]) and 360° scenes: *Playroom* [HPP*18], and *Truck* and *Train* [KPZK17] together with a variety of styles ranging from classical paintings to abstract images: *The Starry Night* by Vincent van Gogh, *The Scream* by Edvard Munch, *The Great Wave off Kanagawa* by Hokusai, *On White II* by Wassily Kandinsky, *The Kiss* by Gustav Klimt, and a colored image of the *Mandelbrot Set* created by Wolfgang Beyer (CC BY-SA). The scenes contain vastly different complexities in the represented stimuli—including highly detailed areas, such as the bookcases in the *Playroom*, but also flat white walls. Our generated results for the forward-facing scenes can be seen in Figure 5 and for the 360° scenes in Figure 6.

5.2. Comparison to the State of the Art

We compare our approach to three state-of-the-art approaches: a recent NeRF-based approach, Artistic Radiance Fields (ARF) [ZKB*22a], a recent zero-shot style transfer method for NeRF scenes, StyleRF [LZX*23a], and a recent pre-print that performs style transfer for Gaussian Splatting scenes, StyleGaussian [LZX*24]. For this comparison, we used their official implementations and pre-trained checkpoints [ZKB*22b, LZX*23b, LZX*24]. We do not compare against Gaussian Splatting in Style [SGC*24] and StylizedGS [ZCY*24] as their implementations are not publicly available at the time of writing this paper.

ARF uses the NNFM loss to transfer artistic style and also utilizes explicit color matching to achieve more accurate colors. This approach



Figure 6: Results generated with our approach, \mathcal{G} -Style, for five 360° scenes (columns) given six style exemplars (rows).

can use different volumetric scene representations as its backbone. However, the only publicly available implementation [LZC*23b] uses TensoRF [CXG*22], which in essence is a 3D voxel-based representation, decomposed into several lower-rank tensors. Other representations showcased in the original paper are neural radiance fields and Plenoxels [FKYT*22], but these could not be tested as they were not publicly available. StyleRF also uses TensoRF as its backbone. This approach is based on embedding high-dimensional features into the structure of the scene and, during rendering, uses them to transfer style. StyleGaussian is similar to StyleRF but uses Gaussian Splats as its backbone.

Qualitative Comparison. A comparison for the forward-facing scenes can be seen in Figure 7 and for the 360° scenes in Figure 8. In these figures, we present the comparison of our method to the available checkpoints of other approaches for different scenes. Our approach exhibits better or, at least, comparable results to the other methods. Both StyleRF and StyleGaussian do not faithfully capture the visual style, because they are not able to synthesize small patterns (as is the case for the *Mandelbrot* style in Figure 7, second and seventh row) nor brushstrokes (for *The Starry Night* and *The Scream* styles) in Figures 7 and 8. These two approaches also fail to recreate any bigger patterns, such as in the two *Truck* scenes of Figure 8, especially in the stylization with *On White*

II. Furthermore, StyleGaussian maintains the original Gaussians as resulting from the reconstruction phase. It is, thus, not able to create any meaningful patterns in undersampled areas, which can be easily seen on the walls in the examples of Figure 7, and on the ground or in the sky in all the examples depicted in Figure 8. Moreover, as shown in the *Train* scene in Figure 8, StyleGaussian produces large colorful Gaussians in the sky which do not match any of the used style images. Thus, for the scenes and styles we chose, StyleRF and StyleGaussian cannot reliably generate style-specific details and patterns; but rather focus on recoloring the scenes. Yet, the visual style of an image goes beyond just the colors.

Similarly to ARF, we produce high-frequency details, but by utilizing the CLIP loss, we can also create bigger patterns. This is something that ARF seems to struggle with. In the case of *The Starry Night* (in both Figure 7 and 8), our stylized scenes contain brush-like patterns, but also moon-like and star-like shapes not present in the results of ARF. When using *The Scream* our method occasionally also creates head-like shapes (see the femur of the *T-Rex* in Figure 7 and the surface of the *Train* in Figure 8). If this is not desired, it could be removed by modifying the style image and/or cropping it. Furthermore, *On White II* consists of simple patterns using only a single color. While our method can capture those patterns and create colorful shapes, ARF only

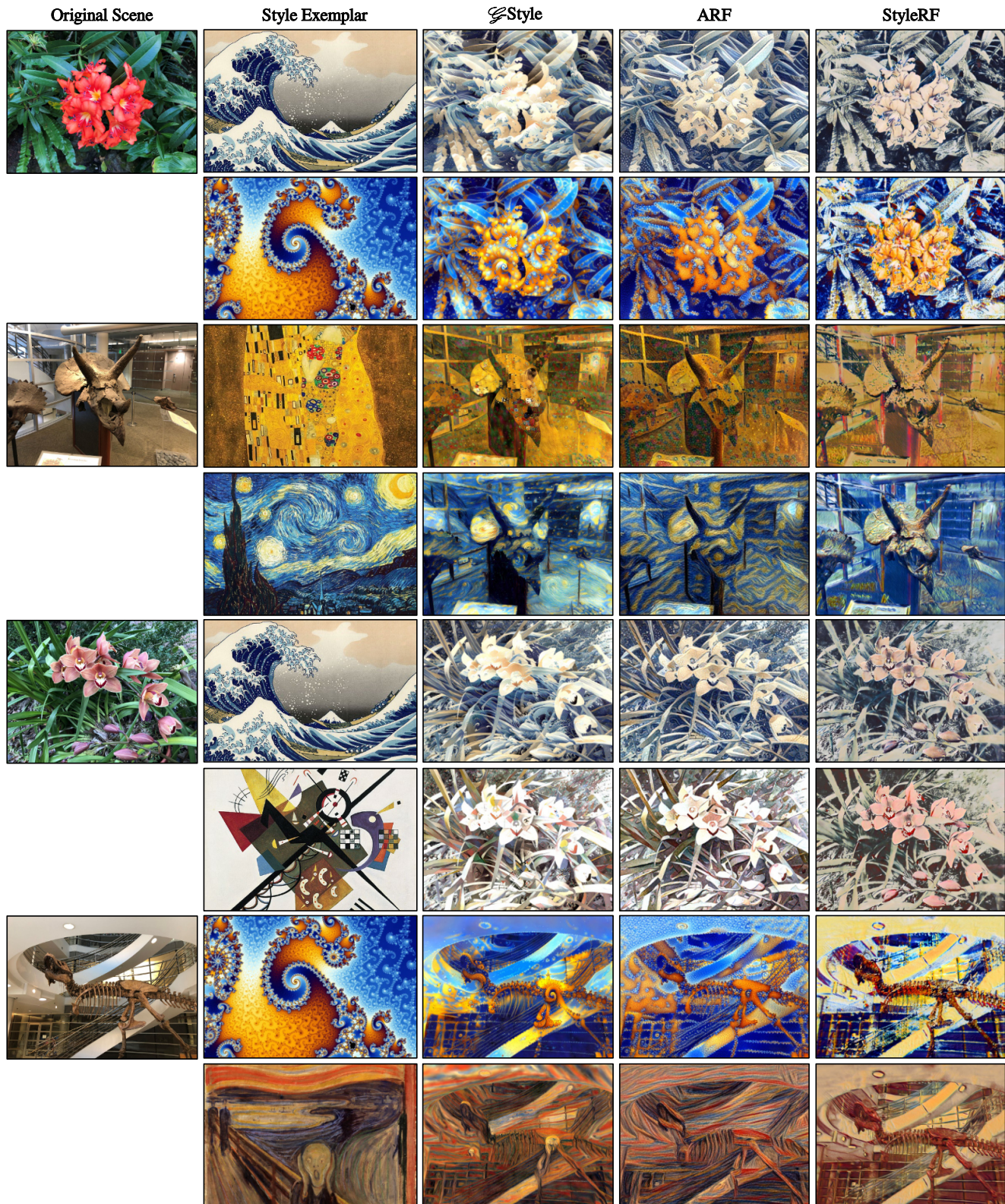


Figure 7: Results generated with our approach (*G*-Style), ARF [ZKB* 22a], and StyleRF [LZC* 23a] (columns) for four forward-facing scenes and two style exemplars for each scene (rows).



Figure 8: Results generated with our approach (*G-Style*), ARF [ZKB*22a], StyleRF [LZC*23a], and StyleGaussian [LZX*24] (columns) for 360° scenes and two style exemplars for each scene (rows).

produces desaturated areas, which are not truthful to the painting. This is evident in the sixth row of Figure 7 and the *Truck* in Figure 8. Finally, when using *The Kiss*, the flower pattern in Figure 7 is less blurry and more recognizable with our method, as opposed to ARF. We conclude that, while ARF can produce highly stylized images that match certain style images, it focuses only on fine details and ignores bigger patterns which may be important to capture a given style successfully.

Speed Comparison. We measured the optimization times of all the compared methods, and we performed all of our measurements on an NVIDIA L4 in Google Colab. The chosen approaches are fundamentally different: our approach and ARF [ZKB*22b] optimize directly the colors of a scene, while StyleRF [LZC*23a] and StyleGaussian [LZX*24] embed VGG features in the scene and later use them during rendering to stylize them. For the tested forward-facing scenes, our pre-processing step takes approximately 5 minutes, and stylization 3–8 minutes depending on the complexity of the scene, the number of ground truth images,

and the size of the style image. For the same scenes, ARF needs 2–7 minutes. For the tested 360° scenes, the preprocessing step of our method takes 8 minutes, and stylization 20–28 minutes. For the same scenes, ARF requires 23–33 minutes. Note that depending on the particular scene and the quality of its reconstruction, we may not need to perform the pre-processing step, as its purpose is to ensure an approximately uniform distribution of the shapes of Gaussians. For StyleGaussian, the process when the embedded features are infused with the style information takes approximately 18 hours per scene, and the rendering of stylized images can be done in real-time. StyleRF needs 30–36 seconds per frame, which includes both style transfer and rendering.

5.3. User Study

To evaluate our method, we conducted an informal, online user study with 24 participants, where we used several of the cases shown in Figures 5–8. We presented each participant with the results of our

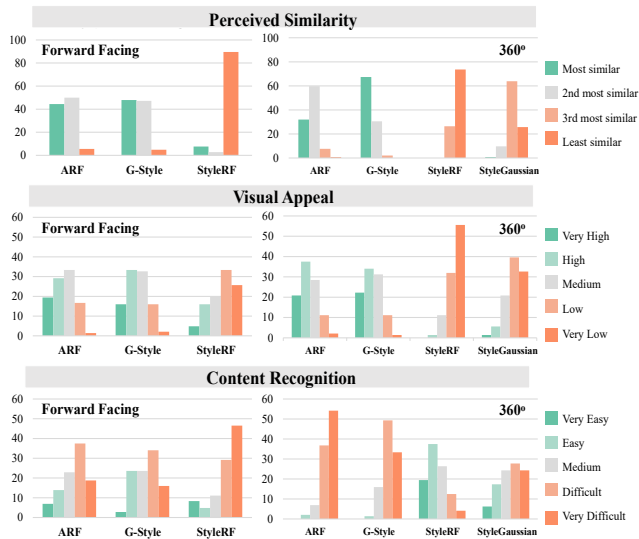


Figure 9: User study results: our approach (\mathcal{G} -Style) vs. ARF [ZKB*22a], StyleRF [LZC*23a], and StyleGaussian [LZX*24].

approach, ARF, StyleRF, and StyleGaussian together with the respective style exemplars and original scenes. Without disclosing any information about any of the approaches, we interviewed the participants to gain feedback about the outputs. Specifically, we asked them to rank the approaches w.r.t. their similarity to the provided style exemplar. For each of the generated results, we also asked the study participants to rate their visual appeal and ability to recognize the original scene content on a 1–5 Likert scale.

The analyzed outcomes of the user study are shown in Figure 9. The study participants ranked our approach as most similar to the style exemplar (47.9% of the participants for the forward-facing scenes vs. 67.4% for the 360° scenes), followed by ARF (44.4% for the forward-facing scenes vs. 32% for the 360° scenes). StyleRF (7.6% for the forward-facing scenes vs. 0% for the 360° scenes) and StyleGaussian (0.7% for the 360° scenes) ranked last. The differences between our approach and ARF are statistically significant only for the 360° scenes, as shown with an ANOVA test followed by pairwise t -tests. The most visually appealing approaches are deemed to be ours and ARF, as indicated visually in the plots of Figure 9. However, the distributions of the ratings of our approach and ARF do not statistically significantly differ. In terms of content recognition, our approach and ARF are comparable to each other. Both are better than StyleRF for forward-facing scenes, but this changes for the 360° scenes, where the ratings of StyleRF (followed by StyleGaussian) are higher. This is to be expected as the effect of the latter two on the stylization of the scenes is less drastic than the former. To sum up, according to our study participants, our approach is comparable to ARF in all investigated aspects—yet, for 360° scenes, ours yields results more faithful to the style.

5.4. Ablations

Our loss does not enforce color transfer and relies on pre-trained networks to synthesize new features. Often, these new features are patterns that do not necessarily have the same color as in the style image. When we do not perform any explicit **color matching**—both during and after

training—synthesized images are discolored (see Figure 10). Conversely, with color matching, the colors are truthful to the style image, as also shown in the work of Zhang et al. [ZKB*22a]. Introducing another style loss, \mathcal{L}_{CLIP} , still does not ensure that the optimization process matches the color statistic in synthesized images. **Without splitting**, the number and shape of Gaussians are the same as after the initial optimization with the ground truth images. This results in blurred areas, where not many Gaussians were originally needed, such as the walls marked in Figure 10. By splitting Gaussians, we are able to introduce details even in those areas. The **content loss** conditions the stylization. This loss helps to keep certain features intact or recognizable, e.g., the letter on the floor in the *Playroom* scene (see Figure 10). Furthermore, it helps to suppress noise, which is not present in the ground truth images. Moreover, the **CLIP loss** enables our method to focus on bigger features. Without it, only very fine features are transferred. For example, without the CLIP loss, the “melting” patterns visible in *The Scream* are not synthesized as depicted in Figure 10. On the other hand, the **NNFM loss** focuses on high-frequency patterns. For instance, the brush strokes are not as visible when switching off the NNFM loss, as depicted in Figure 10. The pre-trained networks used for the style losses were not originally trained for this task but they were trained for classification. As such, using them for this purpose may cause them to produce noise-like artifacts (see Figure 10, no total variation loss). By using the **total variation loss**, we can remove these artifacts. If an overly smooth appearance is desired, this can be achieved by using an even higher weight for this loss, which needs to be tuned for this specific purpose (see Figure 10, high total variation loss).

6. Limitations

By utilizing Gaussian Splatting as the underlying data representation, we also inherit certain visual artifacts that are either caused by their construction process or are fundamental to this representation. Namely, when reconstructing a real-world scene, some “stray” Gaussians may appear as floaters, potentially inhibiting the stylization process. This could be remedied by using a modified version of the original Gaussian Splatting approach [KKLD23], which is currently an active area of research [CW24, FXZ*24, WYZ*24]. Furthermore, splatting can lead to compositing artifacts. Changing the viewpoint slightly may cause a Gaussian to pop in front of another one, given that during rendering the Gaussians are sorted based on the distance from their mean to the camera. Moreover, the stylization process is unguided, which means that an artist lacks full control over the placement of style features. Also, there is no straightforward way to ensure that patterns such as brush strokes point in the desired direction. To the best of our knowledge, there are no reliable ways of selecting style patterns to reconstruct the content of a given image. This seems to be an inherent limitation of style transfer methods based on transferring the distribution of high-level features across images. This could be remedied in future work by considering more advanced style transfer models that can be conditioned to ensure this degree of control.

Furthermore, we rely on pre-trained neural networks to extract features based on which we transfer styles. However, there is no guarantee that the networks can truthfully capture the style features within their latent variables, which could lead to unconvincing results. According to our results shown in Figures 5 and 6, we are able to capture and synthesize the style of varied style images. Yet, with a more robust architecture, we might be able to do so more closely to the original style. Lastly, our approach relies on the stylization of 2D projected areas obtained with



Figure 10: Ablation studies performed for our approach.

a differentiable renderer from a finite set of viewpoints. Depending on the distribution of those viewpoints, certain areas may take longer to converge or might not be fully optimized. Conducting an analysis to identify infrequently viewed areas and strategically placing new cameras in those locations could lead to faster convergence and improved results.

7. Conclusions

We introduced a novel algorithm for stylizing a 3D scene represented by a set of Gaussians to match the style of a given image, \mathcal{G} -Style. By optimizing the geometry of the scene based on the needs of the stylization process and by using a dual loss that captures high and low-frequency style patterns, we generate stylized scenes with higher quality than existing methods in a matter of minutes per style image. In the future, we would like to address the aforementioned limitations, inherited from the Gaussian Splatting representation and the employed neural network architectures, aiming to further improve the quality of the stylized scenes. Finally, we intend to perform additional editing experiments, such as blending multiple styles into one content image or conducting localized or semantic style transfer on distinct regions of the content image, to understand further the strengths and limitations of our proposed approach.

Acknowledgment The authors would like to thank Dr. Bernhard Kerbl for his input at the early stages of this project.

References

- [AHS*21] AN J., HUANG S., SONG Y., DOU D., LIU W., LUO J.: Artflow: Unbiased image style transfer via reversible neural flows. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021). 2
- [CHH24] CHUNG J., HYUN S., HEO J.-P.: Style injection in diffusion: A training-free approach for adapting large-scale diffusion models for style transfer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2024). 2
- [CS16] CHEN T. Q., SCHMIDT M.: Fast patch-based style transfer of arbitrary style. In *arXiv preprint arXiv:1612.04337* (2016). 2
- [CW10] CHEN J., WANG B.: High quality solid texture synthesis using position and index histogram matching. *The Visual Computer* (2010). 2
- [CW24] CHEN G., WANG W.: A survey on 3D gaussian splatting. *arXiv preprint arXiv:2401.03890* (2024). 2, 10
- [CWNN20] CAO X., WANG W., NAGAO K., NAKAMURA R.: PSNet: A Style Transfer Network for Point Cloud Stylization on Geometry and Color. In *2020 IEEE Winter Conference on Applications of Computer Vision (WACV)* (2020). 2
- [CXG*22] CHEN A., XU Z., GEIGER A., YU J., SU H.: Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision (ECCV)* (2022). 7
- [FKYT*22] FRIDOVICH-KEIL S., YU A., TANCIK M., CHEN Q., RECHT B., KANAZAWA A.: Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022). 7
- [FXZ*24] FEI B., XU J., ZHANG R., ZHOU Q., YANG W., HE Y.: 3D gaussian splatting as new era: A survey. *IEEE Transactions on Visualization and Computer Graphics* (2024). 2, 10
- [GLY18] GU S., CHEN C., LIAO J., YUAN L.: Arbitrary style transfer with deep feature reshuffle. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2018). 2
- [GEB15a] GATYS L. A., ECKER A. S., BETHGE M.: A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576* (2015). 1, 2, 5
- [GEB15b] GATYS L. A., ECKER A. S., BETHGE M.: Texture synthesis using convolutional neural networks. *Advances in neural information processing systems* 28 (2015). 2
- [GRGH19] GUTIERREZ J., RABIN J., GALERNE B., HURTUT T.: On Demand Solid Texture Synthesis Using Deep 3D Networks. *Computer Graphics Forum* (2019). 2
- [HB17] HUANG X., BELONGIE S.: Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision* (2017). 2

- [HHY*22] HUANG Y.-H., HE Y., YUAN Y.-J., LAI Y.-K., GAO L.: Stylizednerf: Consistent 3D scene stylization as stylized nerf via 2D-3D mutual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022). 2
- [HJN22] HÖLLEIN L., JOHNSON J., NIESSNER M.: Stylemesh: Style Transfer for Indoor 3D Scene Reconstructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022). 2
- [HMR20] HENZLER P., MITRA N. J., RITSCHEL T.: Learning a Neural 3D Texture Space from 2D Exemplars. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2020). 2
- [HPP*18] HEDMAN P., PHILIP J., PRICE T., FRAHM J.-M., DRETTAKIS G., BROSTOW G.: Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)* (2018). 6
- [HTS*21] HUANG H.-P., TSENG H.-Y., SAINI S., SINGH M., YANG M.-H.: Learning to stylize novel views. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2021). 2
- [JBV17] JETCHEV N., BERGMANN U., VOLLGRAF R.: Texture synthesis with spatial generative adversarial networks. *arXiv preprint arXiv:1611.08207* (2017). 2
- [Ker23] KERBL, BERNHARD AND KOPANAS, GEORGIOS AND LEIMKÜHLER, THOMAS AND DRETTAKIS, GEORGE: 3D Gaussian Splatting for Real-Time Radiance Field Rendering — Implementation. <https://github.com/graphdeco-inria/gaussian-splatting>, 2023. 5
- [KFCO*07] KOPF J., FU C.-W., COHEN-OR D., DEUSSEN O., LISCHINSKI D., WONG T.-T.: Solid Texture Synthesis from 2D Exemplars. In *ACM SIGGRAPH 2007*. ACM, 2007. 2
- [KHR24] KOVÁCS Á. S., HERMOSILLA P., RAIDOU R. G.: Surface-aware mesh texture synthesis with pre-trained 2D cnns. In *Computer Graphics Forum (Eurographics)* (2024). 2
- [KKLD23] KERBL B., KOPANAS G., LEIMKÜHLER T., DRETTAKIS G.: 3D gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics* (2023). 2, 3, 4, 5, 10
- [KPZK17] KNAPITSCH A., PARK J., ZHOU Q.-Y., KOLTUN V.: Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics* (2017). 6
- [KSS19] KOLKIN N., SALAVON J., SHAKHAROVICH G.: Style transfer by relaxed optimal transport and self-similarity. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2019). 2
- [LW16] LI C., WAND M.: Combining markov random fields and convolutional neural networks for image synthesis. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016). 2
- [LYY*17] LIAO J., YAO Y., YUAN L., HUA G., KANG S. B.: Visual attribute transfer through deep image analogy. *ACM Transactions on Graphics* (2017). 2
- [LZC*23a] LIU K., ZHAN F., CHEN Y., ZHANG J., YU Y., SADDIK A. E., LU S., XING E.: StyleRF: Zero-shot 3D style transfer of neural radiance fields. *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2023). 2, 6, 8, 9, 10
- [LZC*23b] LIU K., ZHAN F., CHEN Y., ZHANG J., YU Y., SADDIK A. E., LU S., XING E.: StyleRF: Zero-shot 3D style transfer of neural radiance fields — Implementation. <https://github.com/Kunhao-Liu/StyleRF>, 2023. 6, 7
- [LZL*23] LIU R., ZHAO E., LIU Z., FENG A., EASLEY S. J.: Instant photorealistic style transfer: A lightweight and adaptive approach. *arXiv preprint arXiv:2309.10011* (2023). 2
- [LZX*24] LIU K., ZHAN F., XU M., THEOBALT C., SHAO L., LU S.: StyleGaussian: Instant 3D Style Transfer with Gaussian Splatting. *arXiv preprint arXiv:2403.07807* (2024). 2, 3, 6, 9, 10
- [Max95] MAX N.: Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics* (1995). 3
- [MPSO18] MORDVINTSEV A., PEZZOTTI N., SCHUBERT L., OLAH C.: Differentiable image parameterizations. *Distill* 3, 7 (2018), e12. 2
- [MSOC*19] MILDENHALL B., SRINIVASAN P. P., ORTIZ-CAYON R., KALANTARI N. K., RAMAMOORTHY R., NG R., KAR A.: Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)* (2019). 6
- [MST*20] MILDENHALL B., SRINIVASAN P. P., TANCIK M., BARRON J. T., RAMAMOORTHY R., NG R.: NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *European Conference on Computer Vision (ECCV)* (2020). 1, 2, 3
- [MWWL22] MU F., WANG J., WU Y., LI Y.: 3D photo stylization: Learning to generate stylized novel views from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022). 2
- [NPLX22] NGUYEN-PHUOC T., LIU F., XIAO L.: Snerf: stylized neural implicit representations for 3D scenes. *SIGGRAPH* (2022). 2
- [RKH*21] RADFORD A., KIM J. W., HALLACY C., RAMESH A., GOH G., AGARWAL S., SASTRY G., ASKELL A., MISHKIN P., CLARK J., ET AL.: Learning transferable visual models from natural language supervision. In *International conference on machine learning* (2021). 5
- [SGC*24] SAROHA A., GLADKOVA M., CURRELI C., YENAMANDRA T., CREMERS D.: Gaussian splatting in style. *arXiv preprint arXiv:2403.08498* (2024). 2, 3, 6
- [SZ14] SIMONYAN K., ZISSERMAN A.: Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations (ICLR)* (2014). 5
- [Ull79] ULLMAN S.: The interpretation of structure from motion. *Proceedings of the Royal Society of London. Series B. Biological Sciences* (1979). 3
- [WYZ*24] WU T., YUAN Y.-J., ZHANG L.-X., YANG J., CAO Y.-P., YAN L.-Q., GAO L.: Recent advances in 3D gaussian splatting. *Computational Visual Media* (2024), 1–30. 2, 10
- [WZX23] WANG Z., ZHAO L., XING W.: Stylediffusion: Controllable disentangled style transfer via diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023). 2
- [XCX*24] XU H., CHEN W., XIAO F., SUN B., KANG W.: StyleDyRF: Zero-shot 4D Style Transfer for Dynamic Neural Radiance Fields. *arXiv preprint arXiv:2403.08310* (2024). 2
- [ZCY*24] ZHANG D., CHEN Z., YUAN Y.-J., ZHANG F.-L., HE Z., SHAN S., GAO L.: StylizedGS: Controllable Stylization for 3D Gaussian Splatting. *arXiv preprint arXiv:2404.05220* (2024). 6
- [ZFLS24] ZHANG D., FERNANDEZ-LABRADOR C., SCHROERS C.: CoARF: Controllable 3D Artistic Style Transfer for Radiance Fields. *International Conference on 3D Vision (3DV)* (2024). 3
- [ZGW*22] ZHAO X., GUO J., WANG L., LI F., ZHENG J., YANG B.: STS-GAN: Can We Synthesize Solid Texture with High Fidelity from Arbitrary Exemplars? *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence* (2022). 2
- [ZHT*23] ZHANG Y., HUANG N., TANG F., HUANG H., MA C., DONG W., XU C.: Inversion-based style transfer with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2023). 2
- [ZKB*22a] ZHANG K., KOLKIN N., BI S., LUAN F., XU Z., SHECHTMAN E., SNAVELY N.: ARF: Artistic radiance fields. In *European Conference on Computer Vision* (2022). 2, 5, 6, 8, 9, 10
- [ZKB*22b] ZHANG K., KOLKIN N., BI S., LUAN F., XU Z., SHECHTMAN E., SNAVELY N.: ARF: Artistic radiance fields — Implementation. <https://github.com/Kai-46/ARF-svox2>, 2022. 6, 9