# StylizedGS: Controllable Stylization for 3D Gaussian Splatting

Dingxi Zhang, Yu-Jie Yuan, Zhuoxun Chen, Fang-Lue Zhang, Zhenliang He, Shiguang Shan, and Lin Gao*

**Abstract**—As XR technology continues to advance rapidly, 3D generation and editing are increasingly crucial. Among these, stylization plays a key role in enhancing the appearance of 3D models. By utilizing stylization, users can achieve consistent artistic effects in 3D editing using a single reference style image, making it a user-friendly editing method. However, recent NeRF-based 3D stylization methods encounter efficiency issues that impact the user experience, and their implicit nature limits their ability to accurately transfer geometric pattern styles. Additionally, the ability for artists to apply flexible control over stylized scenes is considered highly desirable to foster an environment conducive to creative exploration. To address the above issues, we introduce StylizedGS, an efficient 3D neural style transfer framework with adaptable control over perceptual factors based on 3D Gaussian Splatting (3DGS) representation. We propose a filter-based refinement to eliminate floaters that affect the stylization effects in the scene reconstruction process. The nearest neighbor-based style loss is introduced to achieve stylization by fine-tuning the geometry and color parameters of 3DGS, while a depth preservation loss with other regularizations is proposed to prevent the tampering of geometry content. Moreover, facilitated by specially designed losses, StylizedGS enables users to control color, stylized scale, and regions during the stylization to possess customization capabilities. Our method achieves high-quality stylization results characterized by faithful brushstrokes and geometric consistency with flexible controls. Extensive experiments across various scenes and styles demonstrate the effectiveness and efficiency of our method concerning both stylization quality and inference speed.

**Index Terms**—Gaussian Splatting, Style Transfer, Perceptual Control

✦

## 1 INTRODUCTION

N OWADAYS, the once professionally-dominated domain of artistic content creation has become increasingly accessible to novice users, thanks to recent groundbreaking advancements in visual artistic stylization research. As a pivotal artistic content generation tool in crafting visually engaging and memorable experiences, 3D scene stylization has attracted growing research efforts. Previous methods have achieved attractive 3D scene style transfer results over diverse explicit representations such as mesh [1], [2], [3], voxel [4], [5], and point cloud [6], [7], [8]. However, the quality of their results is limited by the precision of the geometric reconstructions. The recent 3D stylization methods benefit from the emerging implicit neural representations [9], [10], [11], such as Neural Radiance Fields (NeRF) [12], [13], [14], [15], [16], achieving more faithful and consistent stylization within 3D scenes. Nonetheless, NeRF-based methods are computationally intensive to optimize and suffer from the geometry artifacts of the original radiance fields.

The recently introduced 3D Gaussian Splatting (3DGS) [17], showcasing remarkable 3D reconstruction quality from multi-view images with high efficiency, suggests representing the 3D scene using an array of colored and explicit 3D Gaussians. Performing 3D stylization on scene represented by 3DGS brings the benefits of prompt response and flexible style control to 3D stylization applications due to the explicit nature of 3DGS. Recent 3DGS scene manipulation methods [18], [19], [20], [21] explore the editing and control of 3D Gaussians using text prompts within designated regions of interest or semantic tracing. However, these approaches are driven by text input and fall short of delivering detailed style transfer capabilities. Some styles are also difficult to describe through simple text prompt. Recently appeared studies on 3DGS stylization [22], [23] show limited capabilities of learning and transferring style features. Notably, these works also lack comprehensive control over stylization effects.

In this paper, we introduce the first *controllable* scene stylization method based on 3DGS, StylizedGS. Given a reference style image, our method effectively transfer the style features to the entire 3D scene represented by a set of 3D Gaussians. Our method facilitates the artistic creation of visually coherent novel views that exhibit transformed and detailed style features in a visually reasonable manner. More importantly, StylizedGS operates at a fast inference speed, ensuring efficiency when rendering stylized scenes. To achieve effective stylization and avoid geometry distortion, a straightforward approach is to fine-tune the color of 3D Gaussians, minimizing both style loss and content loss. However, focusing solely on color fails to effectively capture the overall stylistic features when the style patterns contain intricate details. Instead, we propose formulating the 3DGS

---

- * Corresponding Author is Lin Gao (gaolin@ict.ac.cn).
- Dingxi Zhang and Zhuoxun Chen are with the University of Chinese Academy of Sciences, Beijing, China.
  E-Mail: {zhangdingxi20a, zhuoxunchen20}@mails.ucas.ac.cn
- Yu-Jie Yuan and Lin Gao are with the Beijing Key Laboratory of Mobile Computing and Pervasive Device, Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China, and also with the University of Chinese Academy of Sciences, Beijing, China.
  E-Mail: {yuanyujie, gaolin}@ict.ac.cn
- Fang-Lue Zhang is with Victoria University of Wellington, New Zealand.
- Zhenliang He and Shiguang Shan are with the Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China, and also with the University of Chinese Academy of Sciences, Beijing, China.
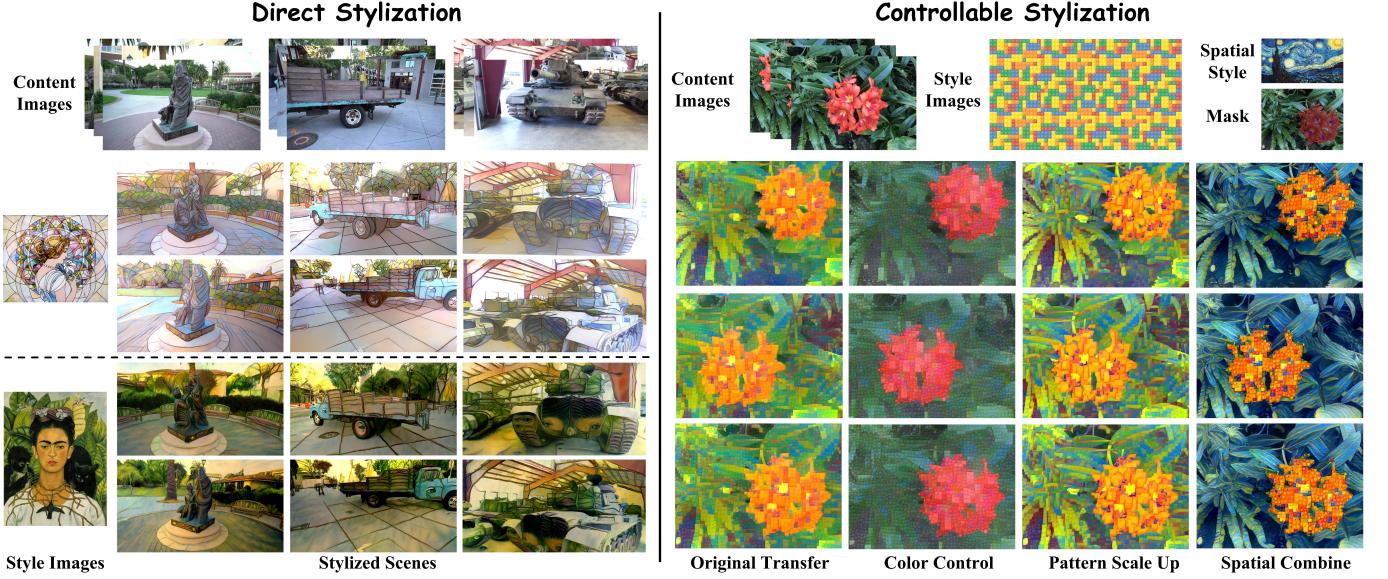  E-mail: {hezhenliang, sgshan}@ict.ac.cn

Fig. 1: **Stylization Results.** Given a 2D style image, the proposed StylizedGS method can stylize the pre-trained 3D Gaussian Splatting to match the desired style with detailed geometric features and satisfactory visual quality within a few minutes. We also enable users to control several perceptual factors, such as color, the style pattern size (scale), and the stylized regions (spatial), during the stylization to enhance the customization capabilities.

stylization as a joint optimization of both the geometry and color of 3D Gaussians to capture the detailed style features while preserving the semantic content. Specifically, our method employs a two-stage stylization framework. In the first stage, we optimize the 3D Gaussians to align the color distribution with the style image while reducing artifacts and geometry noises using our filter-based 3D Gaussian refinement. In the second stage, we achieve final stylization by applying a nearest neighbor feature match (NNFM) loss, optimizing the 3D Gaussians' parameters to capture coherent fine-grained style details.

To enable users to customize the stylization process, we propose a series of strategies and loss functions to control perceptual factors such as color, scale, and spatial regions [24], [25], [26], [27] in the final stylized results. By incorporating this enhanced level of perceptual controllability, our approach synthesizes diverse 3D scene stylization results with varied visual characteristics. Furthermore, we introduce a depth preservation loss that eliminates the necessity for additional networks or regularization operations to preserve the learned 3D scene geometry. It establishes a delicate balance between geometric optimization and depth preservation, facilitating effective style pattern transfer while mitigating significant deterioration of geometric content.

Our contribution can be summarized as follows:

- We introduce StylizedGS, a novel controllable 3D Gaussian stylization method that organically integrates the filter-based 3DGS refinement and the depth preservation loss with other stylization losses to transfer detailed style features and produce faithful novel stylized views.
- We empower users with an efficient stylization process and flexible control through specially designed learning schemes and losses, enhancing their creative

capabilities.
- Our approach achieves significantly reduced training and rendering times while generating high-quality stylized scenes compared with existing 3D stylization methods.

## 2 RELATED WORK

**Image Style Transfer.** Style transfer aims to generate synthetic images with the artistic style of given images while preserving content. Initially proposed in neural style transfer methods [28], [29], this process involves iteratively optimizing the output image using Gram matrix loss and content loss calculated from VGG-Net [30] extracted features. Subsequent works [31], [32], [33], [34] have explored alternative style loss formulations to enhance semantic consistency and capture high-frequency style details such as brushstrokes. Feed-forward transfer methods [35], [36], [37], where neural networks are trained to capture style information from the style image and transfer it to the input image in a single forward pass, ensuring fast stylization. Recent improvements in style loss [13], [33], [34] involve replacing the global Gram matrix with the nearest neighbor feature matrix, improving texture preservation. Some methods adopt patch matching for image generation, like PatchMatch [38] and Fast PatchMatch [39], but can only be applied to limited views. Recent works [40], [41] explore popular tools such as CycleGAN [42], Transformer [43], and diffusion model [44] for style transfer. However, directly applying these 2D methods to 3D stylization will lead to issues like blurriness or view inconsistencies due to the inadequate consideration of 3D geometry.

**3D Gaussian Splatting.** 3D Gaussian Splatting (3DGS) [17] has emerged as an approach for real-time radiance field rendering and 3D scene reconstruction. Recent methods [18],
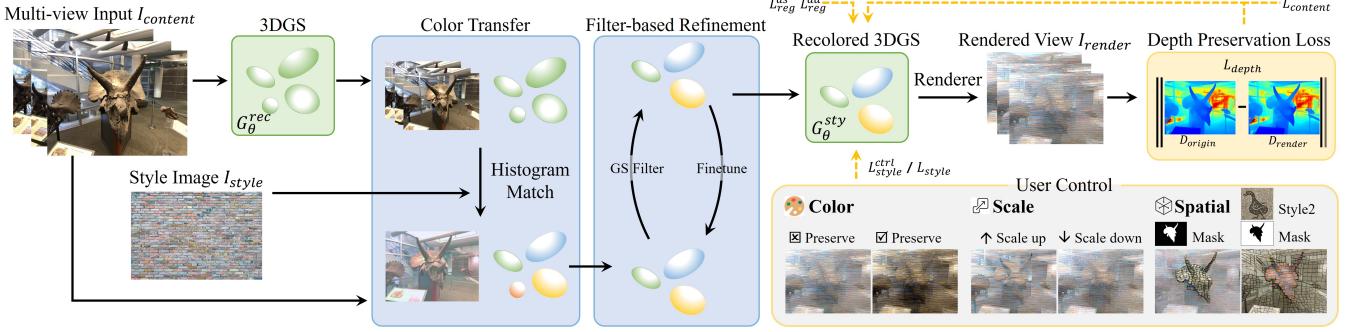
Fig. 2: **StylizedGS Pipeline.** We first reconstruct a photo-realistic 3DGS $G_\theta^{rec}$ from multi-view input. Following this, color matching with the style image is performed, accompanied by the filter-based refinement to preemptively address potential artifacts. During optimization, we employ multiple loss terms to capture detailed local style structures and preserve geometric attributes. Users can flexibly control color, scale, and spatial attributes during stylization through customizable loss terms. Once this stylization is done, we can obtain consistent free-viewpoint stylized renderings.

[19], [20], [45] enhance semantic understanding of 3D scenes and enable efficient text-based editing using pre-trained 2D models. Additionally, geometry deformation [46] and texture editing [47] have also been explored on 3DGS. Please refer to the survey [48], [49] for more details. Despite these advancements, existing 3DGS works lack support for image-based 3D scene stylization that faithfully transfers detailed style features while offering flexible control.

**3D Style Transfer.** With the increasing demand for 3D content, neural style transfer has been expanded to various 3D representations. Stylization on meshes often utilizes differential rendering to transfer style from rendered images to 3D meshes, enabling geometric or texture transfer [1], [2], [3]. Other works, using point clouds as the 3D proxy, ensure 3D consistency when stylizing novel views. LSNV [7] employs featurized 3D point clouds modulated with the style image, followed by a 2D CNN renderer to generate stylized renderings. However, explicit methods' performance is constrained by the quality of geometric reconstructions, often leading to noticeable artifacts in complex real-world scenes.

Implicit methods, such as NeRF [50], have gained considerable attention for their enhanced capacity to represent complex scenes. Many NeRF-based stylization works incorporate image style transfer losses [13], [28] during training or adopt a mutually learned image stylization network [10], [12] to optimize color-related parameters given a reference style image. Approaches like [14], [51] support both appearance and geometric stylization to mimic the reference style, achieving consistent results in novel-view stylization. However, these methods involve time-consuming optimization and exhibit slow rendering due to expensive random sampling in volume rendering. They also lack user-level flexible and accurate perceptual control for stylization.

While controlling perceptual factors, such as color, stroke size, and spatial aspects, have been extensively explored in image style transfer [24], [25], [26], they remain under-developed for 3D stylization. [52], [53] establish semantic correspondence in transferring style across the entire stylized scene, but they primarily focus on spatial control and lack the capability for users to interactively specify arbitrary regions. ARF-plus [27] introduces more perceptual controllability into the stylization of radiance fields, yet the demand for enhanced flexibility and personalized, diverse characteristics in 3D stylization remains unmet.

In this work, leveraging the 3DGS representation, we achieve rapid stylization within a minute of training, ensuring real-time rendering capabilities. Recent arXiv works [22], [23] on 3DGS-based style transfer exhibit limitations in maintaining intricate style details and coherence. [22] is feedforward-based and struggles to capture locally coherent style features and patterns with AdaIN [54], while [23] simply adopts classic style loss to optimize the color of Gaussians. Our method not only effectively captures distinctive details from the style image and preserves recognizable scene content with fidelity, but also empowers users with perceptual control over color, scale, and spatial factors for customized stylization.

## 3 METHOD

Given a 3D scene represented by a collection of captured images with corresponding camera parameters, our objective is to efficiently accomplish consistent style transfer from an arbitrary 2D style image to the 3D scene. Building on recent advances in 3D Gaussian Splatting (3DGS) [17], we optimize the 3DGS representation to generate stylized scenes with intricate artistic features and high visual quality. Our stylization process consists of two key steps: 1) We first recolor the 3D scene to align its color statistics with those of the style image. Simultaneously, we adopt a filter-based 3D Gaussian refinement scheme to minimize the impact of floaters in the reconstruction, which is crucial for ensuring the final stylization quality. 2) We then optimize 3D Gaussians further by exploiting nearest-neighbor feature matching style loss to capture detailed local style patterns. During the optimization, we incorporate a depth preservation loss and regularization terms to preserve the geometric features of the scene. Finally, we introduce how to achieve flexible control over the color, scale, and spatial attributes in our 3DGS stylization (Sec. 3.3). As a result, users can create stylized renderings with customized artistic expression, and explore the consistently stylized scene in a free-view manner. The pipeline of our method is shown in Fig. 2.

### 3.1 Preliminaries: 3D Gaussian Splatting

Gaussian Splatting [17] encapsulates 3D scene information using a set of 3D colored Gaussians. This technique exhibits rapid inference speeds and exceptional reconstruction quality compared to NeRF. To represent the scene, each Gaussian is described by a centroid $p = \{x, y, z\} \in \mathbb{R}^3$, a 3D vector $s \in \mathbb{R}^3$ for scaling, and a quaternion $q \in \mathbb{R}^4$ for rotation. Additionally, an opacity value $\alpha \in \mathbb{R}$ and a color vector $c$ represented in the coefficients of a spherical harmonic (SH) function of degree 3 are used for fast alpha-blending during rendering. These trainable parameters are collectively symbolized by $G_{\theta_i}$, where $G_{\theta_i} = \{p_i, s_i, q_i, \alpha_i, c_i\}$, representing the parameters for the $i$-th Gaussian. To visualize the 3D Gaussians and supervise their optimization, 3DGS projects them onto the 2D image plane. The implementation leverages differentiable rendering and gradient-based optimization on each pixel for the involved Gaussians. The pixel color $c^\alpha$ is determined by blending the colors $c_i$ of those ordered Gaussians that overlap the pixel. This process can be formulated as:

$$c^\alpha = \sum_{i \in N} T_i \alpha_i c_i \tag{1}$$

where $T_i$ is the accumulated transmittance and $\alpha_i$ is the alpha-compositing weight for the $i$-th Gaussian.

### 3.2 Style Transfer to 3D Gaussian Splatting

Given a set of content images $\mathcal{I}_{content}$ that are captured from different viewpoints in the same scene, we first reconstruct its 3DGS model $G_\theta^{rec}$. Our goal is to transform $G_\theta^{rec}$ to a stylized 3DGS model $G_\theta^{sty}$ that matches the detailed style features of a given 2D style image $\mathcal{I}_{style}$ while maintaining the content of the original scene.

**Color Transfer.** We first transfer the color distribution of the style image to the 3DGS $G_\theta^{rec}$, enhancing the alignment of hues with the style image. We will introduce how we enable a flexible color control later. Drawing inspiration from [24], we employ a linear transformation of colors in RGB space, followed by a color histogram matching procedure between the style images and the views of the training set. Let $p_s^i$ represent the set of all pixels in the style image, and $p_c^i$ denote the set of all pixels in the content images to be recolored. We then solve the following linear transformation to align the mean ($\mu_*$) and covariance ($\Sigma_*$) of color distributions between the content image set and the style image:

$$p_c^{re} = \mathbf{A}p_c + b, \; c^{re} = \mathbf{A}c + b$$
$$s.t. \; \mu_{p_c^{re}} = \mu_{p_s}, \Sigma_{p_c^{re}} = \Sigma_{p_s} \tag{2}$$

where $\mathbf{A} \in \mathbb{R}^{3 \times 3}$ and $b \in \mathbb{R}^3$ are the solution to the linear transformation, $p_c^{re}$ is a recolored pixel of the content images. The color parameter $c$ of 3DGS $G_\theta^{rec}$ is transformed to the recolored version $c^{re}$. Please refer to the supplementary document for the mathematical derivation.

**Filter-based Refinement.** After the color transfer step, the color attributes $c^{re}$ of 3DGS and the recolored rendering images $\mathcal{I}_{content}^{re}$ are consistent with the palette of the style image. However, some floaters in the original 3DGS will also be colored during this process, significantly affecting the quality of the stylization, as shown in Fig. 13. These colored floaters often result from Gaussians with excessively

large scales or improper opacity values. To address this issue, we design a filter-based refinement process to exclude these floaters while ensuring reconstruction quality before stylization. We use the recolored content images $\mathcal{I}_{content}^{re}$ as supervision to fine-tune the 3DGS with the recolored color parameter $c^{re}$ and filter the 3D Gaussian floater at every certain number of fine-tuning iterations. The filtering is performed based on the scale and opacity of the 3D Gaussians. Specifically, Gaussians whose sizes are in the top k% or whose opacities are in the lowest k% will be filtered out. The threshold k% remains constant during the fine-tuning process. By maintaining a fixed threshold, we ensure that the reconstructed scene does not change significantly, as the fine-tuning process is relatively short. Such a refinement step minimizes the impact on the overall scene reconstruction quality and enhances color correspondence, while the GS filter removes the floaters. During fine-tuning, we incorporate the reconstruction loss from [17]:

$$\mathcal{L}_{rec} = (1 - \lambda_{rec})\mathcal{L}_1(\mathcal{I}_{content}^{re}, \mathcal{I}_{render}) + \lambda_{rec}\mathcal{L}_{D-SSIM} \tag{3}$$

**Stylization.** We then employ an intuitive and effective strategy to learn 3D Gaussian stylization by leveraging features extracted by a pre-trained convolutional neural network (e.g., VGG [30]). The 3D Gaussians are optimized with a set of loss functions between training views and the style image. To transfer detailed high-frequency style features from a 2D style image to a 3D scene, we exploit the nearest neighbor feature matching concept inspired by [13], [55], where we minimize the cosine distance between the feature map of rendered images and its nearest neighbor in the style feature map. We extract the VGG feature maps $F_{style}$, $F_{content}$, and $F_{render}$ for $\mathcal{I}_{style}$, $\mathcal{I}_{content}^{re}$, and the rendered images $\mathcal{I}_{render}$, respectively. The nearest neighbor feature match loss is formulated as:

$$\mathcal{L}_{style}(F_{render}, F_{style}) = \frac{1}{N}\sum_{i,j}^{N} D(F_{render}(i, j), F_{style}(i*, j*))$$
$$\tag{4}$$

where, $(i*, j*) = \arg\min_{i',j'} D(F_{render}(i, j), F_{style}(i', j'))$, and
$D(a, b)$ is the cosine distance between two feature vectors $a$ and $b$. To preserve the original content structure during the stylization, we additionally minimize the mean squared distance between the content feature $F_{content}$ and the stylized feature $F_{render}$:

$$\mathcal{L}_{content} = \frac{1}{H \times W}||F_{content} - F_{render}||^2 \tag{5}$$

where $H$ and $W$ represent the size of rendered images. While the content loss $\mathcal{L}_{content}$ mitigates the risk of excessive stylization, optimizations applied to the geometric parameters of 3D Gaussians can still induce changes in scene geometry. To address this, we introduce a depth preservation loss to maintain general geometric consistency without relying on any additional depth estimator. It should be noticed that there is a trade-off between geometric optimization and depth preservation loss, allowing our stylization to effectively transfer style patterns while avoiding potential damage to geometric content. Specifically, in the color transfer stage, we initially generate the original depth map $D_{origin}$ by employing the alpha-blending method of
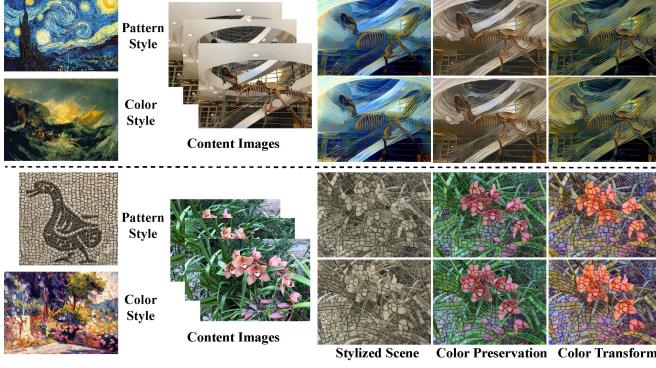
Fig. 3: **Color Control Results.** Our approach facilitates versatile color management in stylized outputs, allowing users to retain the scene's original hues or apply distinct color schemes from alternative style images. Users can choose to transfer the entire style, only the pattern style, or a mix of arbitrary patterns and color styles.



Fig. 4: **Scale Control Results.** Our method enables users to flexibly control the scale of basic style elements, such as adjusting the density of Lego blocks, as demonstrated in the example provided in the last row.

3DGS. The alpha-blended depth $d^\alpha$ of each pixel is calculated as $d^\alpha = \sum_i T_i \alpha_i d_i$, where $d_i$ is the depth value for $i$-th Gaussian. During the stylization, we minimize the $L_2$ loss between the rendered depth map $D_{render}$ and the original one $D_{origin}$:

$$\mathcal{L}_{depth} = \frac{1}{H \times W} ||D_{origin} - D_{render}||^2 \qquad (6)$$

We also add some regularization terms that perform on the changes of scale $\Delta s$ and opacity $\Delta \alpha$ in $G_\theta^{sty}$:

$$\mathcal{L}_{reg}^{ds} = \frac{1}{M}||\Delta s||, \quad \mathcal{L}_{reg}^{d\alpha} = \frac{1}{M}||\Delta \alpha|| \qquad (7)$$

Finally, a total variation loss $\mathcal{L}_{tv}$ is used to smooth the rendered images in the 2D domain. The total loss during the stylization phase is:

$$\begin{aligned} \mathcal{L} = &\lambda_{sty}\mathcal{L}_{style} + \lambda_{con}\mathcal{L}_{content} + \lambda_{dep}\mathcal{L}_{depth} \\ &+ \lambda_{sca}\mathcal{L}_{reg}^{ds} + \lambda_{opa}\mathcal{L}_{reg}^{d\alpha} + \lambda_{tv}\mathcal{L}_{tv} \end{aligned} \qquad (8)$$

where $\lambda_*$ is the corresponding loss weight.

### 3.3 Perceptual Control

Building on the previous components of our approach, we can achieve both rapid and intricate style transfer while preserving the original content. Here, we address the challenge users face in customizing the stylization process. Inspired by [24], we propose a series of strategies and loss functions for 3DGS stylization to control various perceptual factors, including color, scale, and spatial areas, enhancing the customization capabilities of our approach. The stylization loss will be replaced with a different form, namely $L_{style}^{color}$, $L_{style}^{scale}$ or $L_{style}^{spatial}$ defined below.

**Color Control.** Color information within an image is a crucial perceptual aspect of its style. Unlike other stylistic elements such as brush strokes or dominant geometric shapes, color is largely independent and stands out as a distinct characteristic. As illustrated in Fig. 3, users may seek to preserve the original color scheme by transferring only the pattern style or by combining the pattern style of one image with the color style of another. Consequently, our method provides users with independent control over
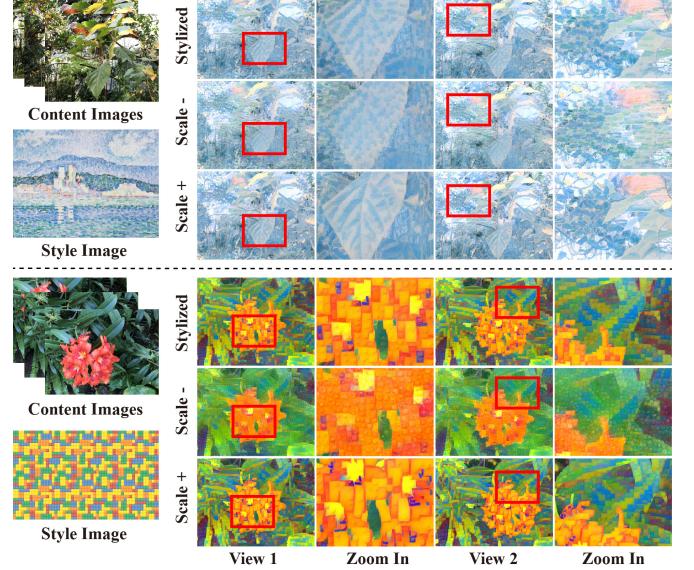
color style during the stylization process. To achieve this, in our color transfer step, we can pre-color 3D scenes to align with any user-desired hues. Similar to the luminance-only algorithm proposed in [56], we utilize the YIQ color space to separate the learning of luminance and color information. As defined in the color style loss $\mathcal{L}_{style}^{color}$ in Eq. 9, we extract the luminance channel from the rendered view $\mathcal{I}_{render}^Y$ and the style image $\mathcal{I}_{style}^Y$ for the luminance-only style loss calculation. Meanwhile, the RGB channels are retained for the content loss $\mathcal{L}_{content}$ calculation.

$$\mathcal{L}_{style}^{color} = \mathcal{L}_{style}(F(\mathcal{I}_{render}^Y), F(\mathcal{I}_{style}^Y)) \qquad (9)$$

**Scale Control.** The scale-related style features, such as the thickness of brushstrokes or the density of the basic grains (as depicted in Fig. 4), are foundational style elements and play a vital role in defining visual aesthetics. Here, the 'scale' we mention refers to the style pattern size. The stroke size, as a typical example of scale features, is influenced by two factors: the receptive field of the convolution in the VGG network and the size of the style image input to the VGG network, as identified by [25]. Our exploration reveals that simply adjusting the receptive field is efficient for achieving scale-controllable stylization in 3DGS. We modulate the size of the receptive field by manipulating the selection of layers in the VGG network and adjusting the weights for each layer during the computation of the scale style loss $\mathcal{L}_{style}^{scale}$. As shown in Eq. 10, the style losses from different layer blocks are combined as $\mathcal{L}_{style}^{scale}$:

$$\mathcal{L}_{style}^{scale} = \sum_{l \in b_s} w_l \cdot \mathcal{L}_{style}(F_l(\mathcal{I}_{render}), F_l(\mathcal{I}_{style})) \qquad (10)$$

Here, $F_l$ represents the feature map at the $l$-th layer within the $s$-th VGG-16 block $b_s$, and $w_l$ controls the corresponding weights at the $l$-th layer.

**Spatial Control.** Traditional 3D style transfer methods typically apply a single style uniformly across the entire scene,
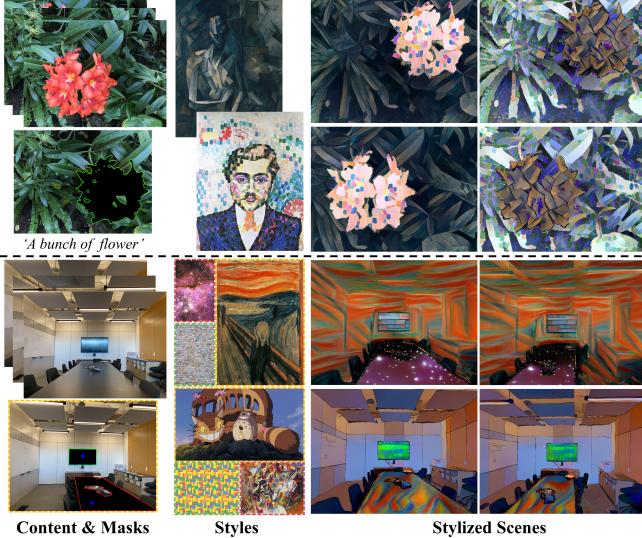
Fig. 5: **Spatial Control Results.** By specifying masks in the content images, users can transfer different styles to desired regions. Users can input text prompt or specify certain points (highlighted in blue) to generate region masks (depicted in black).

resulting in uniform stylization throughout the composition. However, there are scenarios where it becomes necessary to transfer distinct styles to different areas or objects. Users may also wish to specify region-to-region constraints for their preferred style distribution within a scene, as illustrated in Fig.5. We provide point-based interaction using SAM [57] and a proposed mask-tracking strategy, and a language-based approach using the LangSAM method [58] to generate masks across different views. Further details can be found in the supplementary document. Similarly, users can also specify certain areas in the style image to delineate style regions for partial style transfer.

To provide spatial guidance during the stylization process, we introduce a spatial style loss, denoted as $\mathcal{L}_{style}^{spatial}$. Assuming the user specifies $r$ regions in a single view of a scene to be matched with $r$ style regions, the spatial loss $\mathcal{L}_{style}^{spatial}$ is formulated as follows:

$$\mathcal{L}_{style}^{spatial} = \sum_r w_r \cdot \mathcal{L}_{style}(F(M_r^c \circ \mathcal{I}_{render}^r), F(M_r^s \circ \mathcal{I}_{style}^r))$$
(11)

Here, $M_r^c$ and $M_r^s$ represent the binary spatial masks on the rendered views and the $r$-th style image, respectively. Additionally, during the color transfer step, we transform the color information in the region of the style image to the corresponding masked area of the content image to improve color correspondence.

### 3.4 User Interface

We develop an interactive user interface integrated with our proposed method to facilitate convenient image-based scene stylization and various perceptual controls, as illustrated in Fig. 6. Users can upload a style image and multi-view scene images through the right control panel to generate stylized scenes. After the optimization, our interface displays the stylized scenes in the left window which can be viewed in real time. Additionally, the system supports interactive controls and various render modes by adjusting options in the



Fig. 6: **The User Interface of StylizedGS.** Users can upload style and multi-view scene images via the right control panel, and the stylization results are displayed in the left panel in real time after optimization. The interface supports view control, color control, scale control, and spatial control with point-based and language-based interactions.

control panel. Users can change the views by dragging the display window. For spatial control, we have implemented both point-based and language-based interactions to flexibly generate masks. Furthermore, users can customize the stylized scenes by adjusting the optimization steps and loss functions, combining different perceptual controls and style images, thereby supporting detailed and expressive stylization. We also provide a video of manipulating the interactive interface in the supplementary material.

## 4 EXPERIMENTS

We include the implementation details, especially the training settings in our supplementary document.

**Datasets.** We conduct extensive experiments on multiple real-world scenes inLLFF [59] which contains forward-facing scenes and Tanks & Temples (T&T) dataset [60] which includes unbounded 360° large outdoor scenes. Furthermore, our experiments involve a diverse set of style images from [13] and WikiArt [61], allowing us to assess our method's capability to handle a wide range of stylistic exemplars.

**Baselines.** We compare our method with the state-of-the-art 3D stylization methods, including LSNV [7], ARF [13], StyleRF [62], Ref-NPR [53] and StyleGaussian [22], which are based on different 3D representations. Specifically, LSNV is based on point cloud, ARF, StyleRF and Ref-NPR are based on NeRF, and StyleGaussian is based on 3DGS. LSNV, StyleRF and StyleGaussian are feed-forward-based, while others are optimization-based. For all methods, we use their released code and pre-trained models. As Ref-NPR is a reference-based method, we use AdaIN [54] to obtain a stylized reference view according to its paper.

### 4.1 Qualitative Comparisons

We show visual comparisons on the LLFF dataset in Fig. 7. By jointly optimizing the opacity and color components, our method performs better in learning intricate geometric patterns and meso-structures. Specifically, for the scene of "trex" with a Lego style, other methods capture only the color style or a basic square pattern, whereas our approach successfully illustrates the detailed structure of the Lego blocks, clearly depicting even the small cylindrical patterns.
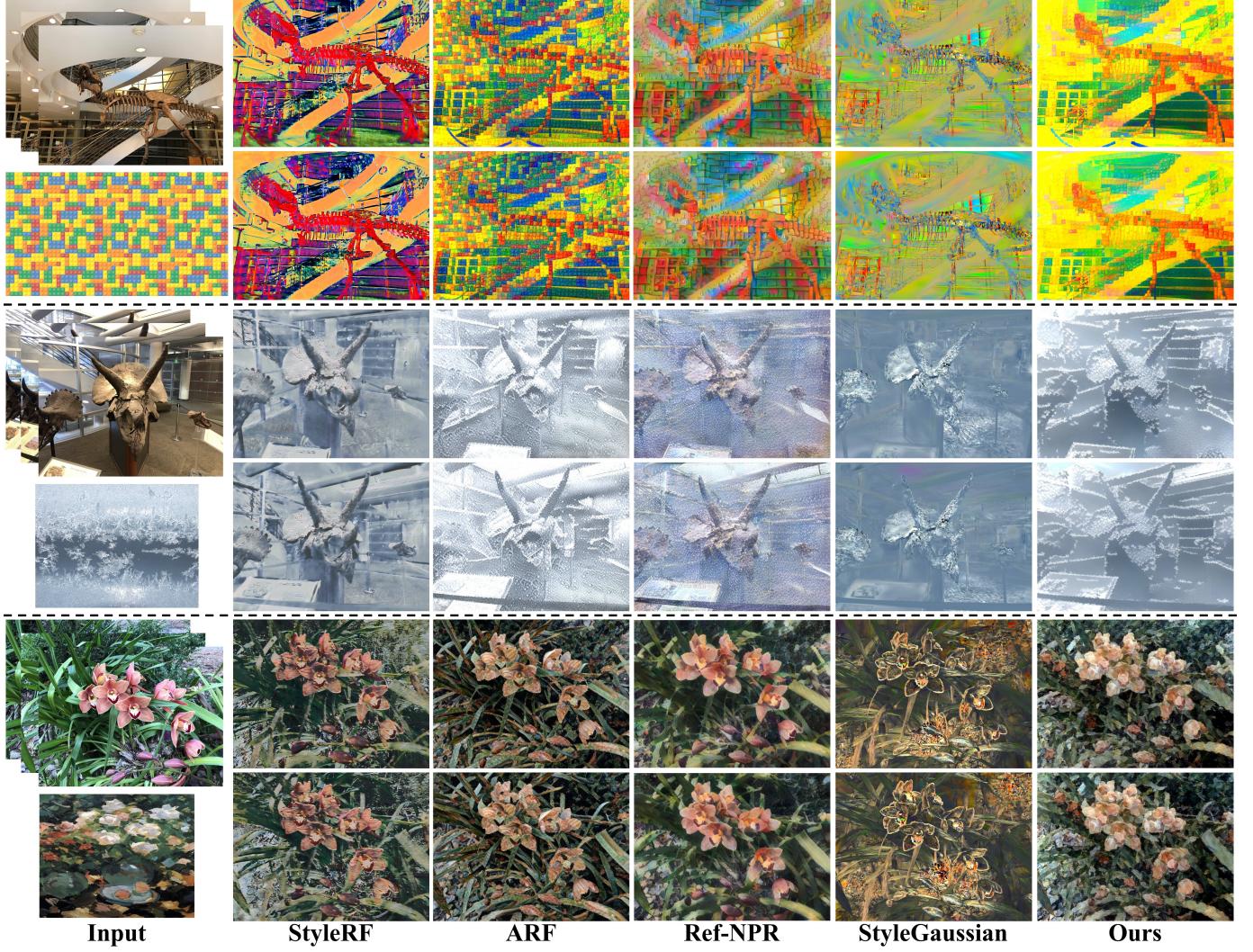
Fig. 7: **Qualitative comparisons with baselines on LLFF dataset.** Compared to other methods, our method excels in learning intricate and accurate geometric patterns while effectively preserving the semantic content of the original scene. (**Please zoom in for better view.**)

This is also evident in the scene of "horns" with a snow style, where our method produces visually more pleasing results that faithfully match the style of the hexagonal snowflake structure. Additionally, our method preserves more semantic information in the stylized scene by utilizing the depth map to maintain spatial distribution in the image. We show more stylization results on the LLFF dataset in Fig. 9.

Fig. 8 presents comparison results on the T&T dataset. It can be seen that our method achieves a better style match with the style image compared to the others. For instance, LSNV tends to produce overly smoothed results and lacks intricate structures, while ARF struggles to consistently transfer the style across the entire scene and exhibits deficiencies in color correspondence, such as the results of the second group. StyleGaussian does not faithfully transfer the style, resulting in blurry results. In contrast, our method accurately captures both the color tones and brushstrokes throughout the entire scene, demonstrating superior style matching compared to the baselines. We show more stylization results on the T&T dataset in Fig. 10.

TABLE 1: Quantitative comparisons on stylization under novel views. We report ArtFID, SSIM, DISTS, average training time (Avg. train), and average rendering FPS (Avg. FPS) for our method and other baselines. '-' denotes that the method does not require individual stylization training.

| Metrics | ArtFID($\downarrow$) | SSIM($\uparrow$) | DISTS($\downarrow$) | Avg. train($\downarrow$) | Avg. FPS($\uparrow$) |
|---|---|---|---|---|---|
| LSNV | 52.75 | 0.13 | 0.33 | - | 0.71 |
| StyleRF | 40.61 | 0.41 | 0.30 | - | 5.8 |
| ARF | 35.73 | 0.29 | 0.31 | 2.25 min | 8.2 |
| Ref-NPR | 33.56 | 0.31 | 0.33 | 2.54 min | 7.3 |
| StyleGaussian | 37.92 | 0.35 | 0.30 | - | 142 |
| Ours | **28.29** | **0.55** | **0.26** | **0.87 min** | **153** |

## 4.2 Quantitative Comparisons

A significant advantage of the 3DGS representation is the capability of rendering high-quality views in real-time and with much less training time. Therefore, we conducted quantitative comparisons with other methods to demonstrate both the high fidelity and efficiency of our approach. As the optimization time during the stylization process is notably influenced by the resolution, we randomly selected
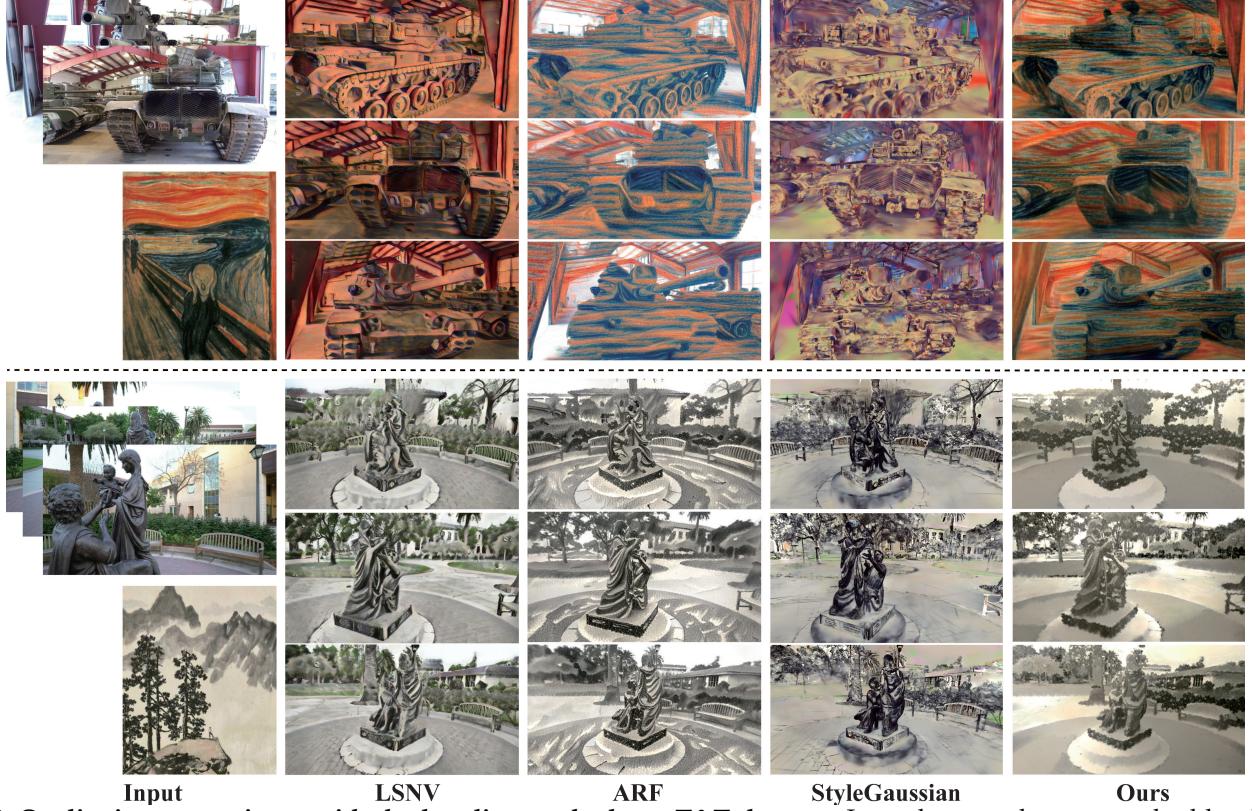
Fig. 8: **Qualitative comparisons with the baseline methods on T&T datasets.** It can be seen that our method has better results which faithfully capture both the color styles and pattern styles across different views. (**Please zoom in for better view.**)
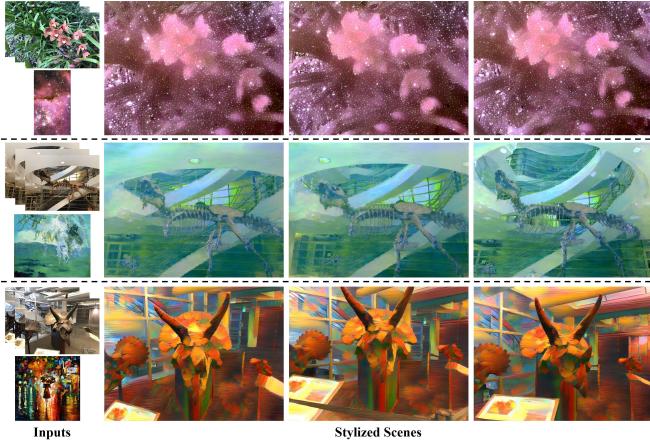


Fig. 9: **Stylization results on LLFF dataset.** Our method demonstrates robust performance, effectively adapting to various styles and diverse scenes.



Fig. 10: **Stylization results on T&T dataset.** It can be seen that our method faithfully captures both the color styles and pattern styles across different views.

50 stylized scenes with consistent resolution from the LLFF dataset and evaluated the average times. For the stylization metrics, we choose ArtFID [63] to measure the quality of stylization following [64]. We also use SSIM [65] and DISTS [66] between the content and stylized images to measure the performance of detail preservation and structural similarity following [35], [67], [68], [69], [70]. As shown in Tab. 1, our method excels in both stylized visual quality and content preservation. Regarding efficiency, our method has significant advantages in average rendering FPS (Avg. FPS)
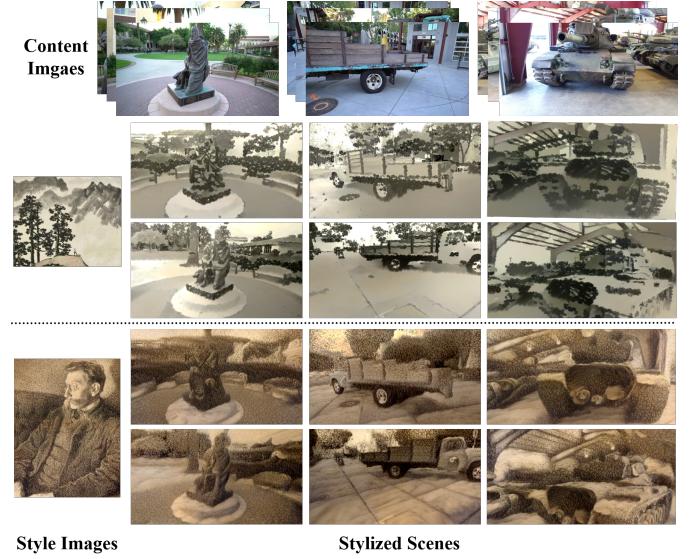
and can achieve real-time free-view synthesis. At the same time, it outperforms other optimization-based stylization methods in terms of average training time (Avg. train) for stylization and can achieve per-scene stylization within 1 minute on a single NVIDIA RTX 4090 GPU. Note that while StyleRF and StyleGaussian are feed-forward methods and

do not require individual stylization training for each case, it's worth mentioning that StyleRF's rendering is slow and cannot be viewed in real time.
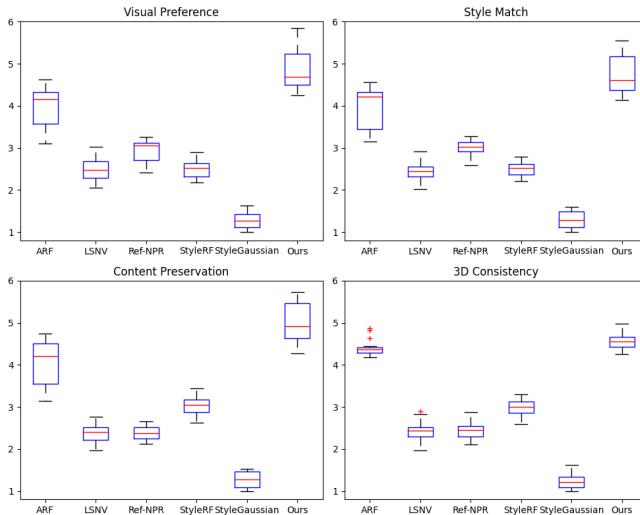


Fig. 11: **User Study.** We record the user preference in the form of a boxplot. Our results obtain more preferences in visual preference, style match level, content preservation level, and 3D consistency quality than other state-of-the-art stylization methods.

**User Study** We conducted a user study comparing our method with other baseline methods. We randomly selected 20 sets of stylized views of 3D scenes from both the LLFF and T&T datasets, processed by various methods, and invited 25 participants. Participants were then asked to vote on the video based on four evaluation indicators: visual preference, style match level, content preservation level, and 3D consistency quality. More details are included in the supplementary document. We collected a total of 1000 votes for all evaluation indicators and presented the results in Fig. 11 in the form of boxplots. Our method, StylizedGS, outperformed other methods across all evaluation indicators, indicating visually preferable and consistent results. A p-value of 0.00001 demonstrates the statistical significance of the obtained results.

### 4.3 Controllable Stylization Results

Our method enables users to control stylization in three ways, including color, scale, and spatial region. The corresponding controllable stylization results are shown in Figs. 3, 4 and 5, respectively. Fig. 12 shows an example where a sequential series of controls are applied to a scene across different conceptual factors. In addition to the direct stylization result, we sequentially apply different controls from left to right: preserving the color of the original scene, increasing the scale of the pattern style, and adopting spatial control to apply two styles. Such flexibility allows for controllable refinement of key parameters or stylistic attributes, satisfying specific aesthetic preferences or artistic requirements and fostering enhanced artistic expression.

### 4.4 Ablation Study

We conducted ablation studies to validate our design choices. The original scene and input style image are presented in Fig. 13 (a). When the recoloring procedure is not

TABLE 2: Quantitative ablation study results. We report ArtFID and SSIM for our method and other ablations over 50 randomly chosen stylized cases.

| Metrics | w/o recolor | w/o density | w/o depth | w/o GS filter | Ours |
|---|---|---|---|---|---|
| ArtFID (↓) | 31.57 | 38.54 | 31.20 | 30.02 | **28.29** |
| SSIM (↑) | 0.54 | 0.53 | 0.35 | 0.39 | **0.55** |

applied, as seen in (b), the stylized scene exhibits inferior color correspondence with the style image. In Fig. 13 (d)(e), we illustrate the importance of our depth preservation loss in preserving the original scene's geometry. Without applying depth loss, the background of the Horse scene disappears, and the Trex's skeleton becomes blurrier. The heatmap in the top left corner depicts the difference between the rendered depth maps and ground truth depth maps. The notable differences in the heatmap of (d) further emphasize the efficacy of our depth preservation loss. In Fig. 13 (c), we highlight the importance of fine-tuning both the color and density components in 3DGS, which is crucial for learning intricate style patterns. Neglecting the optimization of the density makes it challenging to capture stroke details and texture patterns in the style image. Therefore, a balance is established wherein fine-tuning the density component harmonizes with the depth preservation loss, thereby facilitating the transfer of intricate style details while keeping the integrity of the original geometry.

Additionally, we investigate how the proposed 3DGS filter affects the reconstructed scenes. Normally, floaters become more noticeable after the style transfer process. As shown in Fig. 14, the texture of stones and leaves extend into the air around the objects without the 3DGS filter. In contrast, the results of our method exhibit enhanced clarity, with fewer visible colored floaters. Tab. 2 shows the quantitative results of the ablation studies, demonstrating the contributions of all components of our method.

## 5 DISCUSSION AND CONCLUSION

**Limitations.** Although our method achieves efficient and controllable 3DGS stylization, it still has some limitations. First, the geometric artifacts from the original 3DGS reconstruction may impact the quality of the final stylized scenes, as shown in Fig. 15. Although our filter-based refinement can eliminate some floaters, it cannot eliminate all geometric artifacts. In the future, we will incorporate some improved 3DGS reconstruction methods that add geometric constraints, such as SuGaR [71], to address this issue. In addition, our optimization-based method cannot achieve instant stylization, However, this does not compromise the effectiveness of our method or its practical value, as the optimization process typically takes only about 1 minute.

In this paper, we introduce StylizedGS, the first controllable 3D scene stylization method based on the 3D Gaussian Splatting representation. Following 3DGS reconstruction, our proposed filter-based refinement minimizes the impact of floaters in the reconstruction, which is crucial for achieving the desired stylization effect. Additionally, we propose the adoption of nearest-neighbor feature matching style loss to optimize both geometry and color parameters of 3D Gaussians, enabling the capture of detailed style features and facilitating 3D scene stylization. We further
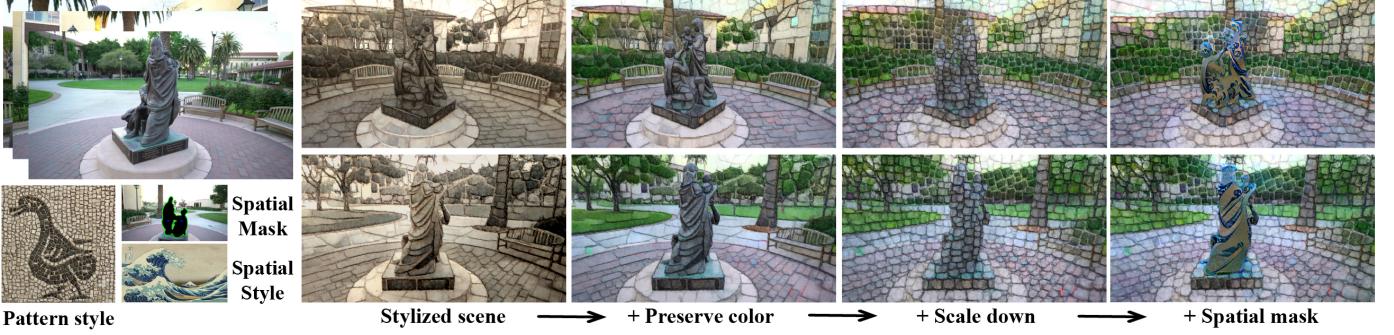
Fig. 12: **Sequential control results.** Given multiple control conditions, we can achieve a sequence of controllable stylization. We first show a stylization result and then, from left to right, progressively implement controls that preserve the color of the original scene, increase the scale of the pattern style, and adopt spatial control to apply two styles.
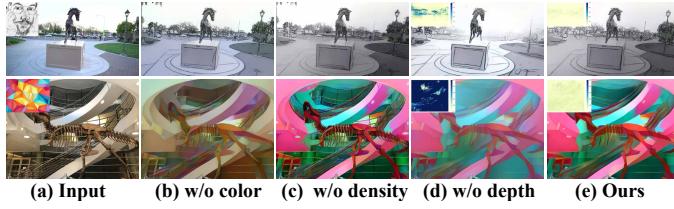


Fig. 13: **Ablation Study about density and color control.** 'w/o density' presents the stylized scene without density component fine-tuning in 3DGS. 'w/o recolor' displays the stylized results without applying recoloring while 'w/o depth' shows the results without incorporating depth loss, leading to the disappearance of original geometry and semantic content.
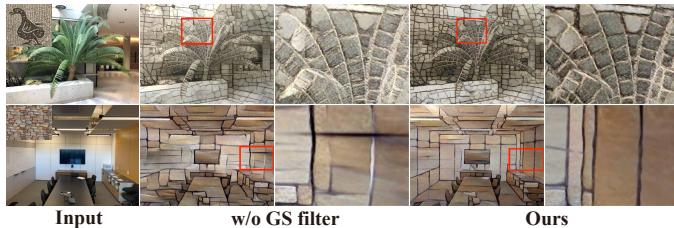


Fig. 14: **Ablation study about 3D Gaussian filter.** The filter effectively helps eliminate artifacts.

introduce a depth preservation loss for regularization to maintain the overall structure during stylization. Moreover, we design controllable methods for three perceptual factors: color, stylization scale, and spatial regions, providing users with specific and diverse control options. Qualitative and quantitative experiments demonstrate that our method outperforms existing 3D stylization methods in terms of effectiveness and efficiency.



Fig. 15: **Limitation**. The geometric artifacts from the original 3DGS reconstruction may impact the quality of the final stylized scenes.
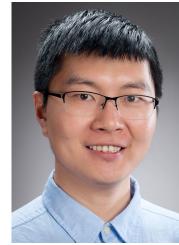
# REFERENCES

[1] H. Kato, Y. Ushiku, and T. Harada, "Neural 3d mesh renderer," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[2] O. Michel, R. Bar-On, R. Liu, S. Benaim, and R. Hanocka, "Text2mesh: Text-driven neural stylization for meshes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 13 492–13 502.

[3] K. Yin, J. Gao, M. Shugrina, S. Khamis, and S. Fidler, "3dstylenet: Creating 3d shapes with geometric and texture style variations," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12 456–12 465.

[4] J. Guo, M. Li, Z. Zong, Y. Liu, J. He, Y. Guo, and L.-Q. Yan, "Volumetric appearance stylization with stylizing kernel prediction network." *ACM Trans. Graph.*, vol. 40, no. 4, pp. 162–1, 2021.

[5] O. Klehm, I. Ihrke, H.-P. Seidel, and E. Eisemann, "Property and lighting manipulations for static volume stylization using a painting metaphor," *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 7, pp. 983–995, 2014.

[6] X. Cao, W. Wang, K. Nagao, and R. Nakamura, "Psnet: A style transfer network for point cloud stylization on geometry and color," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer vision*, 2020, pp. 3337–3345.

[7] H.-P. Huang, H.-Y. Tseng, S. Saini, M. Singh, and M.-H. Yang, "Learning to stylize novel views," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 13 869–13 878.

[8] E. Bae, J. Kim, and S. Lee, "Point cloud-based free viewpoint artistic style transfer," in *2023 IEEE International Conference on Multimedia and Expo Workshops (ICMEW)*. IEEE, 2023, pp. 302–307.

[9] P.-Z. Chiang, M.-S. Tsai, H.-Y. Tseng, W.-S. Lai, and W.-C. Chiu, "Stylizing 3d scene via implicit representation and hypernetwork," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 1475–1484.

[10] T. Nguyen-Phuoc, F. Liu, and L. Xiao, "Snerf: stylized neural implicit representations for 3d scenes," *ACM Transactions on Graphics (TOG)*, vol. 41, no. 4, pp. 1–11, 2022.

[11] Z. Fan, Y. Jiang, P. Wang, X. Gong, D. Xu, and Z. Wang, "Unified implicit neural stylization," in *European Conference on Computer Vision*. Springer, 2022, pp. 636–654.

[12] Y.-H. Huang, Y. He, Y.-J. Yuan, Y.-K. Lai, and L. Gao, "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*, 2022, pp. 18 342–18 352.

[13] K. Zhang, N. Kolkin, S. Bi, F. Luan, Z. Xu, E. Shechtman, and N. Snavely, "Arf: Artistic radiance fields," in *ECCV*, 2022.

[14] C. Wang, R. Jiang, M. Chai, M. He, D. Chen, and J. Liao, "Nerf-art: Text-driven neural radiance fields stylization," *IEEE Transactions on Visualization and Computer Graphics*, 2023.

[15] H.-W. Pang, B.-S. Hua, and S.-K. Yeung, "Locally stylized neural radiance fields," in *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE Computer Society, 2023, pp. 307–316.

[16] Z. Zhang, Y. Liu, C. Han, Y. Pan, T. Guo, and T. Yao, "Transforming

radiance field with lipschitz network for photorealistic 3d scene stylization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 20712–20721.

[17] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics*, vol. 42, no. 4, July 2023. [Online]. Available: https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/

[18] M. Ye, M. Danelljan, F. Yu, and L. Ke, "Gaussian grouping: Segment and edit anything in 3d scenes," *arXiv preprint arXiv:2312.00732*, 2023.

[19] J. Fang, J. Wang, X. Zhang, L. Xie, and Q. Tian, "Gaussianeditor: Editing 3d gaussians delicately with text instructions," *arXiv preprint arXiv:2311.16037*, 2023.

[20] Y. Chen, Z. Chen, C. Zhang, F. Wang, X. Yang, Y. Wang, Z. Cai, L. Yang, H. Liu, and G. Lin, "Gaussianeditor: Swift and controllable 3d editing with gaussian splatting," 2023.

[21] J. Tang, J. Ren, H. Zhou, Z. Liu, and G. Zeng, "Dreamgaussian: Generative gaussian splatting for efficient 3d content creation," *arXiv preprint arXiv:2309.16653*, 2023.

[22] K. Liu, F. Zhan, M. Xu, C. Theobalt, L. Shao, and S. Lu, "Stylegaussian: Instant 3d style transfer with gaussian splatting," *arXiv preprint arXiv:2403.07807*, 2024.

[23] A. Saroha, M. Gladkova, C. Curreli, T. Yenamandra, and D. Cremers, "Gaussian splatting in style," *arXiv preprint arXiv:2403.08498*, 2024.

[24] L. A. Gatys, A. S. Ecker, M. Bethge, A. Hertzmann, and E. Shechtman, "Controlling perceptual factors in neural style transfer," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3985–3993.

[25] Y. Jing, Y. Liu, Y. Yang, Z. Feng, Y. Yu, D. Tao, and M. Song, "Stroke controllable fast style transfer with adaptive receptive fields," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 238–254.

[26] C. Castillo, S. De, X. Han, B. Singh, A. K. Yadav, and T. Goldstein, "Son of zorn's lemma: Targeted style transfer using instance-aware semantic segmentation," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 1348–1352.

[27] W. Li, T. Wu, F. Zhong, and C. Oztireli, "Arf-plus: Controlling perceptual factors in artistic radiance fields for 3d scene stylization," *arXiv preprint arXiv:2308.12452*, 2023.

[28] L. A. Gatys, A. S. Ecker, and M. Bethge, "A neural algorithm of artistic style," *arXiv preprint arXiv:1508.06576*, 2015.

[29] ——, "Image style transfer using convolutional neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2414–2423.

[30] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[31] E. Risser, P. Wilmot, and C. Barnes, "Stable and controllable neural texture synthesis and style transfer using histogram losses," *arXiv preprint arXiv:1701.08893*, 2017.

[32] S. Gu, C. Chen, J. Liao, and L. Yuan, "Arbitrary style transfer with deep feature reshuffle," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8222–8231.

[33] N. Kolkin, J. Salavon, and G. Shakhnarovich, "Style transfer by relaxed optimal transport and self-similarity," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10051–10060.

[34] J. Liao, Y. Yao, L. Yuan, G. Hua, and S. B. Kang, "Visual attribute transfer through deep image analogy," *arXiv preprint arXiv:1705.01088*, 2017.

[35] J. An, S. Huang, Y. Song, D. Dou, W. Liu, and J. Luo, "Artflow: Unbiased image style transfer via reversible neural flows," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 862–871.

[36] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1501–1510.

[37] D. Y. Park and K. H. Lee, "Arbitrary style transfer with style-attentional networks," in *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 5880–5888.

[38] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman, "Patchmatch: A randomized correspondence algorithm for structural image editing," *ACM Trans. Graph.*, vol. 28, no. 3, p. 24, 2009.

[39] T. Q. Chen and M. Schmidt, "Fast patch-based style transfer of arbitrary style," *arXiv preprint arXiv:1612.04337*, 2016.

[40] Y. Zhang, N. Huang, F. Tang, H. Huang, C. Ma, W. Dong, and C. Xu, "Inversion-based style transfer with diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 10146–10156.

[41] Y. He, L. Chen, Y.-J. Yuan, S.-Y. Chen, and L. Gao, "Multi-level patch transformer for style transfer with single reference image," in *International Conference on Computational Visual Media*. Springer, 2024, pp. 221–239.

[42] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.

[43] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.

[44] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10684–10695.

[45] J. Cen, J. Fang, C. Yang, L. Xie, X. Zhang, W. Shen, and Q. Tian, "Segment any 3d gaussians," *arXiv preprint arXiv:2312.00860*, 2023.

[46] L. Gao, J. Yang, B.-T. Zhang, J.-M. Sun, Y.-J. Yuan, H. Fu, and Y.-K. Lai, "Mesh-based gaussian splatting for real-time large-scale deformation," *arXiv preprint arXiv:2402.04796*, 2024.

[47] T. Wu, J.-M. Sun, Y.-K. Lai, Y. Ma, L. Kobbelt, and L. Gao, "Deferredgs: Decoupled and editable gaussian splatting with deferred shading," *arXiv preprint arXiv:2404.09412*, 2024.

[48] T. Wu, Y.-J. Yuan, L.-X. Zhang, J. Yang, Y.-P. Cao, L.-Q. Yan, and L. Gao, "Recent advances in 3d gaussian splatting," *arXiv preprint arXiv:2403.11134*, 2024.

[49] G. Chen and W. Wang, "A survey on 3d gaussian splatting," *arXiv preprint arXiv:2401.03890*, 2024.

[50] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[51] S. Xu, L. Li, L. Shen, and Z. Lian, "Desrf: Deformable stylized radiance field," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 709–718.

[52] M. Kumar, N. Panse, and D. Lahiri, "S2rf: Semantically stylized radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2952–2957.

[53] Y. Zhang, Z. He, J. Xing, X. Yao, and J. Jia, "Ref-npr: Reference-based non-photorealistic radiance fields for controllable scene stylization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 4242–4251.

[54] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4401–4410.

[55] N. Kolkin, M. Kucera, S. Paris, D. Sykora, E. Shechtman, and G. Shakhnarovich, "Neural neighbor style transfer," *arXiv e-prints*, pp. arXiv–2203, 2022.

[56] A. H. C. J. N. Oliver, B. Curless, and D. Salesin, "Image analogies," in *Proc. SIGGRAPH 2001*, 2001, pp. 327–340.

[57] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick, "Segment anything," *arXiv:2304.02643*, 2023.

[58] "Language segment-anything," 2023, https://github.com/luca-medeiros/lang-segment-anything.

[59] B. Mildenhall, P. P. Srinivasan, R. Ortiz-Cayon, N. K. Kalantari, R. Ramamoorthi, R. Ng, and A. Kar, "Local light field fusion: Practical view synthesis with prescriptive sampling guidelines," *ACM Transactions on Graphics (TOG)*, 2019.

[60] A. Knapitsch, J. Park, Q.-Y. Zhou, and V. Koltun, "Tanks and temples: Benchmarking large-scale scene reconstruction," *ACM Transactions on Graphics (ToG)*, vol. 36, no. 4, pp. 1–13, 2017.

[61] W. R. Tan, C. S. Chan, H. Aguirre, and K. Tanaka, "Improved artgan for conditional synthesis of natural image and artwork," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 394–409, 2019. [Online]. Available: https://doi.org/10.1109/TIP.2018.2866698

[62] K. Liu, F. Zhan, Y. Chen, J. Zhang, Y. Yu, A. El Saddik, S. Lu, and E. P. Xing, "Stylerf: Zero-shot 3d style transfer of neural radiance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 8338–8348.

[63] M. Wright and B. Ommer, "Artfid: Quantitative evaluation of neural style transfer," in *DAGM German Conference on Pattern Recognition*. Springer, 2022, pp. 560–576.

[64] S. Huang, J. An, D. Wei, J. Luo, and H. Pfister, "Quantart: Quantizing image style transfer towards high visual fidelity," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5947–5956.

[65] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.

[66] K. Ding, K. Ma, S. Wang, and E. P. Simoncelli, "Image quality assessment: Unifying structure and texture similarity," *CoRR*, vol. abs/2004.07728, 2020. [Online]. Available: https://arxiv.org/abs/2004.07728

[67] K. Hong, S. Jeon, H. Yang, J. Fu, and H. Byun, "Domain-aware universal style transfer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14609–14617.

[68] J. Yang, F. Guo, S. Chen, J. Li, and J. Yang, "Industrial style transfer with large-scale geometric warping and content preservation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 7834–7843.

[69] S. Kim, S. Kim, and S. Kim, "Deep translation prior: Test-time training for photorealistic style transfer," *arXiv preprint arXiv:2112.06150*, 2021.

[70] Y. Yu, D. Li, B. Li, and N. Li, "Multi-style image generation based on semantic image," *The Visual Computer*, vol. 40, pp. 1–16, 08 2023.

[71] A. Gu'edon and V. Lepetit, "Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering," *ArXiv*, vol. abs/2311.12775, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:265308825

**Fang-Lue Zhang** is currently a senior lecturer with Victoria University of Wellington, New Zealand. He received the Doctoral degree from Tsinghua University in 2015. His research interests include image and video editing, computer vision, and computer graphics. He received Victoria Early-Career Research Excellence Award in 2019 and Fast-Start Marsden Grant from New Zealand Royal Society in 2020. He is on the editorial board of Computer & Graphics. He is a committee member of IEEE Central New Zealand Section.



**Zhenliang He** is currently an Assistant Professor in the Visual Information Processing and Learning research group at the Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS). He received the Ph.D. from ICT, CAS in 2021 and received the bachelor's degree from Beijing University of Posts and Telecommunications in 2015. His current research focuses on generative models and representation learning.



**Dingxi Zhang** received the bachelor's degree in Computer Science and Technology from University of Chinese Academy of Sciences in 2024. She is currently pursuing a Master's degree at Department of Computer Science, ETH Zurich. Her research interests include 3D vision and computer graphics.



**Shiguang Shan** received M.S. degree in computer science from the Harbin Institute of Technology, Harbin, China, in 1999, and Ph.D. degree in computer science from the Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS), Beijing, China, in 2004. He joined ICT, CAS in 2002 and became a full Professor in 2010. He is now the director of the Key Lab of Intelligent Information Processing of CAS.



**Yu-Jie Yuan** received the bachelor's degree in mathematics from Xi'an Jiaotong University in 2018. He is currently a Ph.D. candidate in the Institute of Computing Technology, Chinese Academy of Sciences. His research interests include computer graphics and 3D vision.



**Lin Gao** received the bachelor's degree in mathematics from Sichuan University and the PhD degree in computer science from Tsinghua University. He is currently an Associate Professor at the Institute of Computing Technology, Chinese Academy of Sciences. He has been awarded the Newton Advanced Fellowship from the Royal Society and the Asia Graphics Association young researcher award. His research interests include computer graphics and geometric processing.



**Zhuoxun Chen** received a bachelor's degree in Computer Science and Technology from the University of Chinese Academy of Sciences in 2024. He is currently pursuing a Master's degree at the Institute of Computing Technology, Chinese Academy of Sciences. His research interests include neural rendering and robotics.