

StrokeFaceNeRF: Stroke-based Facial Appearance Editing in Neural Radiance Field

Xiao-Juan Li^{1,2} Dingxi Zhang^{2,1} Shu-Yu Chen¹ Feng-Lin Liu^{1,2*}

¹Beijing Key Laboratory of Mobile Computing and Pervasive Device,
 Institute of Computing Technology, Chinese Academy of Sciences

²University of Chinese Academy of Sciences

{lixiaojuan, chenshuyu, liufenglin21s}@ict.ac.cn zhangdingxi20a@mails.ucas.ac.cn



Figure 1. StrokeFaceNeRF enables versatile high-quality color stroke-based appearance editing for facial NeRFs. Users can draw strokes over a rendered facial image from any angle and define the editing region using masks. The third column showcases the photo-realistic free-view rendering results with exceptional color correspondence.

Abstract

Current 3D-aware facial NeRF generation approaches control the facial appearance by text, lighting conditions or reference images, limiting precise manipulation of local facial regions and interactivity. Color stroke, a user-friendly and effective tool to depict appearance, is challenging to edit 3D faces because of the lack of texture, coarse geometry representation and detailed editing operations. To solve the above problems, we introduce StrokeFaceNeRF, a novel stroke-based method for editing facial NeRF appearance. In order to infer the missing texture and 3D geometry information, 2D edited stroke maps are firstly encoded into the EG3D’s latent space, followed by a transformer-based editing module to achieve effective appearance changes while preserving the original geometry in editing regions. Notably, we design a novel geometry loss function to ensure surface density remains consistent during training. To further enhance the local manipulation accuracy, we propose a stereo fusion approach which lifts the 2D mask (inferred from strokes or drawn by users) into 3D mask volume, al-

lowing explicit blending of the original and edited faces. Extensive experiments validate that the proposed method outperforms existing 2D and 3D methods in both editing reality and geometry retention.

1. Introduction

Realistic 3D human face generation is a popular topic and has wide range applications in character design, virtual meeting and education. Traditional 3D facial modeling methods rely on software like Maya and Zbrush to construct mesh models while employing texture mapping for facial appearances. However, these methods are time-consuming and labor-intensive, especially when editing local appearance details and aiming for photorealistic results. Neural Radiance Field (NeRF) [23], a new 3D representation technique, easily render photo-realistic images with the input of multi-view images. The further combination of Generative Adversarial Networks (GAN) [10] with NeRF [5, 11, 24, 30] has enabled the generation of high-quality facial NeRF through random sampling. Nevertheless, precise control over appearance details remains a challenge.

*Corresponding Author

Efforts to enhance controllability involve disentangling geometry and appearance during facial NeRF generation. Sun et al. introduced the FE-NeRF [32] that employs separate latent codes to control the intricacies of facial geometry and appearance. Subsequent advancements [16, 31] utilize decoupling techniques based on triplane representations, leading to the facial generation results with significantly enhanced quality. Despite their success in good geometry and appearance swapping, these methods heavily rely on random sampling or reference images for global appearance editing. This reliance poses challenges in the obtainment of desired reference images and precludes precise control over localized appearance attributes, such as selective hair color alteration.

Compared to reference images, colored strokes, functioning as an interactive tool, offer a precise depiction of local appearance details. In stroke-based image editing, one category of methods [17, 26] relies on conditional GAN and completion framework to edit local regions with the stroke guidance. However, artifacts may arise in the editing boundary and result in less realistic images for large-area editing operations, as shown in Figure 6. Another category of methods [22, 27] employs pre-trained facial generation models as decoders, mapping colored strokes to the latent space of the generation model to generate edited facial results. However, these methods struggle to preserve the unchanged geometric features of the input face during editing, as demonstrated in Figure 6. Importantly, these techniques are specifically tailored for 2D facial image editing and cannot be directly applied to 3D facial radiance fields.

This paper presents StrokeFaceNeRF, a framework for editing facial NeRF appearance based on colored strokes, enabling precise modifications of color and lighting in 3D faces. Leveraging 2D facial image generation techniques such as StyleGAN2 [21] and semantic segmentation methods [38], we construct a dataset of faces with identical geometry but varying appearances, generating corresponding colored stroke maps using bilateral filtering and semantic clustering algorithms. To achieve appearance editing, directly encoding edited stroke maps into the latent space of the pre-trained facial NeRF generation model [5] leads to texture smoothness and undesirable geometry distortion due to high-frequency information absence and user’s inaccurate drawing of strokes. To solve this problem, we propose a transformer-based editing network, which fuses the original and stroke-encoded latent codes to modify appearance while preserving the original geometry in edited regions. Furthermore, we introduce a geometry loss term during training to ensure surface density remains unchanged, improving detail preservation. To further maintain unchanged features in non-edited areas for precise editing, we design a volumetric fusion algorithm that integrates localized regions before and after editing in 3D space, enabling local-

ized appearance editing effects in facial models.

Our contributions are summarized as follows:

- We propose the first color stroke-based facial NeRF editing framework to generate free-view realistic appearance editing facial images.
- We design a novel geometry density loss to train the transformer-based latent editing network, which effectively preserves the geometry detail during facial appearance modifications.
- We introduce a volumetric fusion algorithm to precisely edit local regions while maintaining the 3D consistency in unedited regions.

2. Related Work

2.1. Facial Image Editing with Colored Strokes

Existing 2D facial image editing methods using colored strokes can be divided into two categories. One group of methods utilizes conditional Generative Adversarial Networks (cGANs) to accomplish stroke-guided image generation. Sangkloy et al. [29] utilize sketches and colored strokes to directly generate realistic images through a translation network. However, this method fails to achieve localized area editing. To solve this problem, FaceShop [26] utilizes image completion techniques and 2D mask marking for localized facial image editing based on sketches and colored strokes. Jo et al. [18] further introduced SC-FEGAN that employs gated convolutions and new style losses for training, enabling a broader range of editing operations. Xiao et al. [36] design a hair generation network with self-attention mechanisms, achieving more natural and realistic hair editing effects. However, these methods are sensitive to input conditions and result in smooth edits lacking texture information in edited regions. Additionally, they tend to produce artifacts at the editing boundary due to the image completion techniques.

Another category of methods utilizes pre-trained generative models to generate high quality faces instead of pixel-wise translation. Bau et al. [3] first utilize reverse optimization to achieve image shape and color editing in the GAN’s latent space. pSp [27] further designs an image encoding network to map various types of input images, such as sketches and color strokes, into StyleGAN2’s latent space to generate corresponding facial images. In order to generate highly editable latent codes, Tov et al. [34] further constrained the latent codes with a discriminator. However, these methods, as they represent the entire image with latent vectors, struggle to achieve localized editing effects. Instead of GAN dependence, Meng et al. [22] employed pre-trained diffusion models to achieve facial image editing based on colored stroke maps. Images with user-drawn strokes are added Gaussian noise, then progressively denoised to synthesize realistic facial images. However, it is

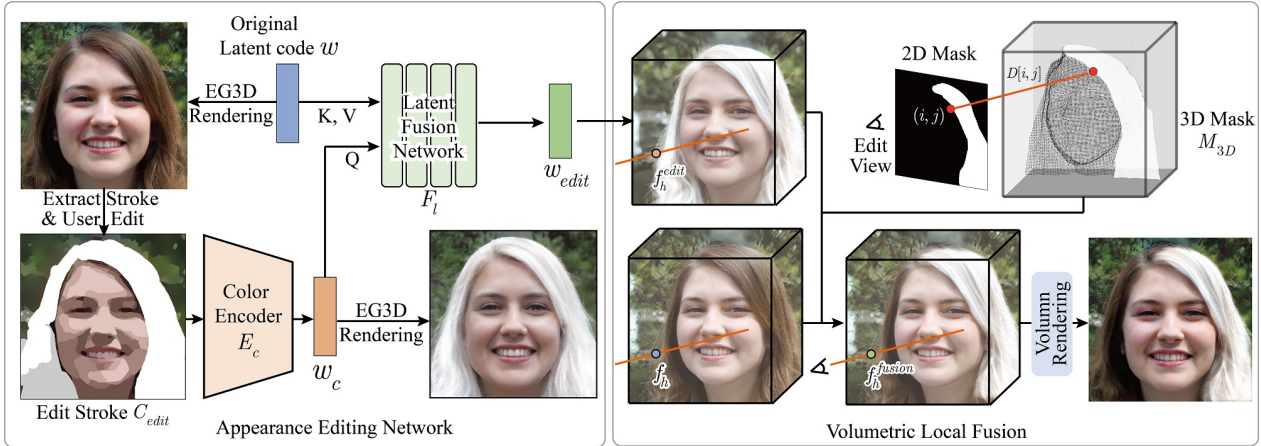


Figure 2. 3D Facial Appearance Editing Framework. Given the original latent code w , free-view images are rendered by EG3D. Users introduce editing operations, resulting in the creation of a color stroke image C_{edit} . This stroke map is encoded by Color Encoder E_c to generate w_c , which is fused with original latent code w by latent fusion network to generate w_{edit} . Utilizing an additional 2D mask, the volumetric fusion approach facilitates targeted appearance editing within specific regions.

challenging to directly apply this method to facial neural radiance field editing. Additionally, because of the direct input of coarse stroke maps, it is challenging for the above methods to maintain geometric features while performing appearance modification in editing regions.

2.2. Editing Facial NeRF

Combining generative adversarial networks with neural radiance fields enables the synthesis of realistic facial images from arbitrary viewpoints. Graf [30] first adds extra latent codes to NeRF’s fully connected network to support face generation and use multiscale discriminators for supervision. Giraffe [24] designs a new framework that employs volume rendering to obtain feature maps and a convolutional framework to generate facial images. StyleSDF [25] and GRAM [7] utilize representation methods based on SDF and implicit planes, respectively, producing higher-quality and more 3D consistent results. EG3D [5] propose a facial NeRF generation framework based on triplane representation and StyleGAN2 [21] generator. However, these methods can only synthesize faces through random sampling and lack the ability for editing.

To generate controllable facial NeRF, Jo et al. [17] use sketches, black-and-white images, and text as control conditions. Gao et al. [9] propose a facial NeRF generation and editing framework based on sketches to enhance synthesis quality and achieve fine editing. However, this method can only edit facial geometry. Another line of work decouples facial NeRF’s geometry and appearance using semantic segmentation. FE-NeRF [32] achieves geometric editing through reverse latent code optimization and appearance editing through random sampling. IDE-3D [31] constructs geometry triplanes and appearance triplanes, and designs a semantic map encoder to generate higher-quality facial results. NeRFFaceEditing [16] apply AdaIN [14] to triplanes

features, enabling the disentanglement of geometry and appearance control. Although the above methods achieve effective appearance modification, reference images or latent codes are required instead of user’s flexible interaction. More importantly, the above methods are designed for the global appearance editing and cannot be directly utilized for local editing.

3. Methodology

We propose a framework, StrokeFaceNeRF, enables users to utilize colored strokes for facial NeRF appearance editing. We employ the EG3D [5] generator architecture, a proficient 3D-aware Generative Adversarial Network (GAN) that combines the StyleGAN2 [21] architecture and neural rendering to produce high-quality 3D shapes. Given a latent code, a triplane representation can be obtained from a StyleGAN2 backbone. The neural renderer then aggregates features from these triplanes and predicts feature images based on a specified camera pose. Subsequently, a super-resolution module is applied to upsample and refine these raw neurally rendered images, resulting in high-quality, multiview-consistent images. More details of the generator can be found in [5].

Our method comprises three key components: dataset construction, colored stroke editing network, and 3D local fusion module. Initially, we introduce a dataset construction approach to acquire facial images with consistent geometry but varying appearances, along with their corresponding colored stroke images as detailed in Section 3.1. Subsequently, an appearance editing network, described in Section 3.2, is established within the latent space of EG3D to predict the 3D faces. Additionally, in Section 3.3, we propose a 3D local fusion method to enable effective appearance editing while meticulously preserving features in

non-edited regions.

3.1. Dataset Construction

To supervise the training of network models, the training data needs to simulate the process of facial appearance editing based on colored strokes. Thus, the training tuples include original faces, edited faces, and corresponding colored stroke maps. Dataset construction involves three steps as shown in Figure 3: (1) building a collection of facial latent codes in EG3D’s latent space, obtained through random sampling and real images inversion; (2) constructing facial tuples with identical geometry but different appearances. Two different methods were employed: one based on semantic segmentation and the other based on StyleGAN2 generative networks; (3) generating colored stroke maps using bilateral filtering and semantic clustering methods.

We utilize two approaches to synthesize paired images with the same geometry but different appearance, indicated as the edited ground truth faces. *Semantic Segmentation-Based Augmentation.* Given the original rendered facial images, we utilize BiSeNet [38] to get semantic segmentation masks. Then, for each local semantic region, we represent it into HSV space, where the Hue channel is uniformly replaced with random color, while the Saturation and Value channels are added random perturbation. Finally, these regions are merged into original images to imitate appearance editing. *StyleGAN2-Based Augmentation.* We further utilize pre-trained generator to synthesize more diverse and realistic training images. Given the original rendered images, ReStyle [2] encodes them into StyleGAN2 latent space, followed with style mixing to change the facial appearance while preserving the same geometry. More details can be found in Supplementary.

With the paired appearance modification images, corresponding color strokes are required to imitate real editing situation. Typically, users drawn strokes are simple and monochromatic. Therefore, synthesized color stroke images must preserve the original image’s color and geometric layout while eliminating high-frequency texture features and local details. To achieve this objective, we employ the bilateral filtering method [33] to effectively preserves the edge features of the facial image while eliminating high-frequency details. Specially, the initial image undergoes median filtering, followed by iterative application of bilateral filtering, as explained in the Supplementary. Nevertheless, some regions still exhibit smooth transitions between different colored areas after filtering, posing challenges for users aiming to replicate specific effects. To tackle this issue, we employ a pixel value clustering technique based on semantic information. Initially, the facial image undergoes semantic segmentation. Within the hair area, clustering allows all pixel colors to be replaced by the nearest cluster center color, thereby eliminating smooth transition regions.

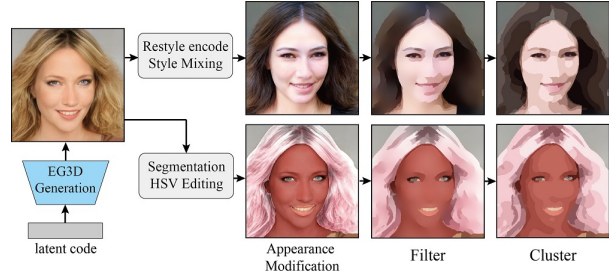


Figure 3. Dataset construction pipeline. Given EG3D rendering images, the appearance is changed by two different approaches. Then, the color stroke images are synthesized by image filter and cluster.

The same methodology is applied to the facial area, yielding the final color stroke image.

3.2. Appearance Editing Network

In our framework, the target face NeRF to be edited is represented as a latent code w within EG3D’s latent space. The latent code can be obtained either by random sampling in EG3D or by applying PTI [28] to invert real images. Utilizing StyleGAN generator and volume rendering, free view facial images can be generated, followed by the approaches in Section 3.1 to extract color stroke images. Users then modify these color stroke images to obtain the edited color stroke image C_{edit} . As illustrated in Figure 2, the input to the appearance editing network consists of the original latent variable w , the edited color stroke image C_{edit} , and the output is the edited latent variable w_{edit} .

The appearance editing network comprises two main components: the **Color Encoder** E_c and the **Latent Fusion Network** F_l . Specifically, given the edited color stroke image C_{edit} , the color encoder that has a similar structure to [27], predicts the latent code w_c . While this latent code generates facial results with color editing effects, it fails to preserve the original facial features, such as the geometric structure of the hair and the original skin color in non-edited areas, as shown in Figure 2. Therefore, a latent code fusion network is designed to utilize the information from the original latent variable w , thus better preserving the original facial features. The latent fusion network F_l adopts a cross-attention transformer block to predict the fused latent variable w_{edit} . This design implicitly extracts the color information from w_c and geometry information from w , guaranteeing a thorough integration of both color-edited and original facial features. In this configuration, the latent code w_c is employed as query tokens, whereas the original latent code w serves as key and value tokens. We also normalize over queries instead of keys as described in [13] to ensure a continuous refinement of w_c during the fusion process with value tokens while mitigating the potential influence of w being too large.

$$w_{edit} = F_l(w, E_c(C_{edit})) \quad (1)$$

3.3. Volumetric Local Fusion

While direct rendering by w_{edit} yields good editing effects within targeted editing areas, as depicted in Figure 7, noticeable undesirable discrepancies emerge in unedited regions. Moreover, intricate editing tasks, such as selectively dyeing specific portions of hair or drawing specific patterns on faces, pose challenges for the latent code-based representation because of its overarching global control inherited from StyleGAN2 backbone. To address these issues, we propose a volumetric fusion approach to enable more refined editing operations.

In order to facilitate local editing, our method introduces support for an additional 2D mask input that indicates the editing region. This 2D mask is further lifted into 3D mask volume to support volumetric fusion. Specifically, we propose a depth guidance approach to sample a set of points near the facial surface. In the volume rendering, for each ray $r_{i,j}$ of the pixel at the i th row and j th column, the camera’s position and ray direction are denoted as $o_{i,j}$ and $d_{i,j}$, respectively. After rendering the original faces, we get the corresponding depth value $D[i, j]$. Then, for each ray, we uniformly sample K points with offset Δx_k within range $[-0.1, 0.1]$. The sampled 3D points set S for the original face is denoted as:

$$S = \cup_{i,j} \cup_{k=1}^K (o_{i,j} + (D[i, j] + \Delta x_k) \cdot d_{i,j}) \quad (2)$$

where the default value of K is 5. Notably, as illustrated in Figure 2, we selectively employ the rays corresponding to pixels within the 2D mask regions.

Subsequently, this generated 3D point set undergoes a transformation into a 3D volumetric mask. Specifically, we initialize a $256 \times 256 \times 256$ 3D volume with zero values throughout the space. The coordinates of the 3D points, along with their neighboring positions, are designated as 1, indicating 3D editing regions. To ensure smoother boundaries, we utilize a three-dimensional dilation and Gaussian filtering, yielding the 3D mask denoted as M_{3D} . Additionally, We implement a mask padding strategy, which introduces direction offsets to the camera poses and extends the mask boundary. This prevents incomplete editing artifacts during view changes. Further details regarding the padding technique are available in the Supplemental material.

To render realistic editing images, the original latent variable w and the edited latent variable w_{edit} are used to synthesize the original and edited triplane features, respectively. For the sampling points h in space during volume rendering, the corresponding original feature f_h and edited feature $f_{\text{edit}}(h)$ are generated. To achieve the editing fusion effect, the features are fused based on the 3D mask M_{3D} :

$$f_{\text{fusion}}(h) = f_{\text{edit}}(h) \cdot M_{3D}[h] + f(h) \cdot (1 - M_{3D}[h]) \quad (3)$$

The fused features f_{fusion} are passed through the decoder of EG3D to obtain color and density information. Afterward,

they are processed through the volume rendering and up-sampling modules to generate the final edited faces.

3.4. Training Strategy

During the training of Appearance Editing Network, the geometry details should be preserved while the appearance editing effectiveness should not be influenced. This is non-trivial because the 2D supervision lacks detailed 3D structure constrain and the synthesized training data has minor geometry distortion because of Style-mixing’s limitation. To solve this problem, we train our network with the follow strategy.

Geometry Density Loss. Based on our NeRF representation, we design a novel density loss to enforce the 3D geometry consistency between the original and edited faces. Specifically, leveraging the 3D points generation approach described in Equation 2, a collection of 3D points is generated. Notably, we use all the rays during the rendering instead of rays in specific regions. These 3D sample points are further projected onto triplane space to query corresponding triplane features that are subsequently decoded into color and density through volume rendering. We minimize the density value between the original and edited faces to effectively preserve the facial geometry while decoupling the appearance:

$$L_{\text{density}} = \sum_{x \in S, x' \in S'} \|\psi_{\text{density}}(\mathcal{F}(x)) - \psi_{\text{density}}(\mathcal{F}'(x'))\|_1 \quad (4)$$

where \mathcal{F} denotes the triplane query process and ψ_{density} denotes density decoder. x' and S' are the 3D points and set for edited faces. These points are sampled with the same offset but new depth map D_{edit} of edited faces. Edited triplane features \mathcal{F}' synthesized by w_{edit} are used for density calculation.

We also utilizes the same loss functions as the pSp framework [27] to match the original stroke distribution. The loss function includes the per-pixel L_2 loss between the ground truth image I_{GT} (obtained in Sec.3.1) and the edited image generated by EG3D G :

$$L_{\text{img}} = \|I_{GT} - G(w_{\text{edit}})\|_2 \quad (5)$$

Here, $\|\cdot\|_2$ denotes the Euclidean distance. To better constrain the reconstruction quality, a perceptual distance loss term is further employed:

$$L_{LPIPS} = \|F(I_{GT}) - F(G(w_{\text{edit}}))\|_2 \quad (6)$$

Where F represents an image feature extraction method based on VGG [19]. In order to better constrain the images generated to preserve the identity features of the original faces, an identity constraint term is employed:

$$L_{ID} = 1 - \langle R(I_{GT}), R(G(w_{\text{edit}})) \rangle \quad (7)$$

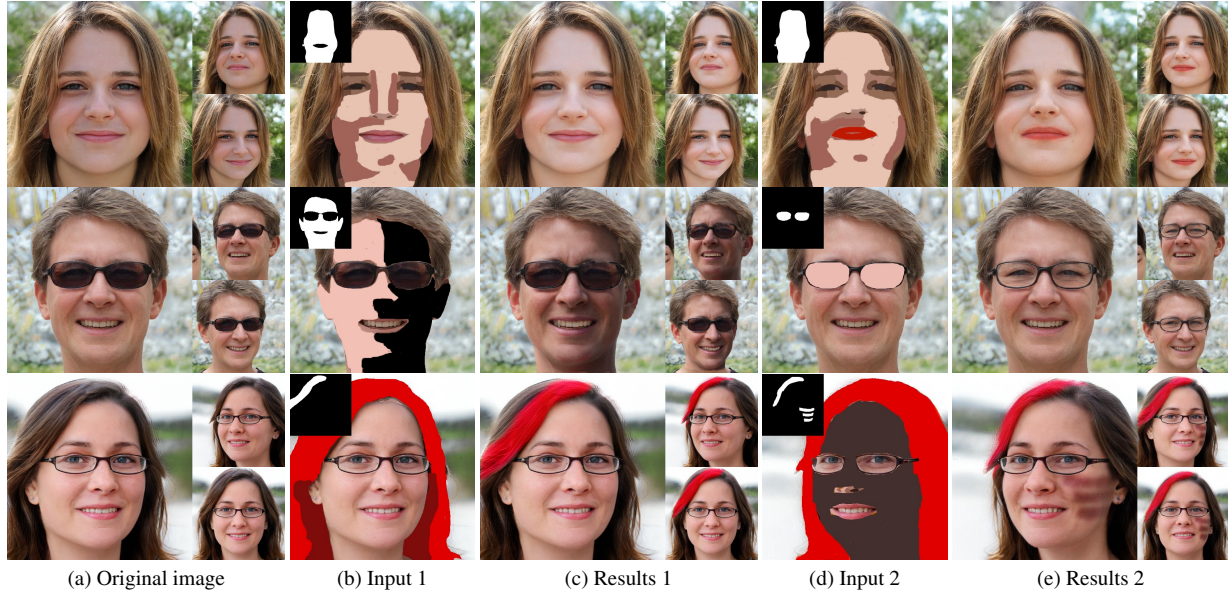


Figure 4. Color stroke-based facial NeRF appearance editing results. Given the original faces (a), users drawn the editing color strokes (b) (d), our methods generates the corresponding editing results (c)(e). Additional masks are shown in top left corner, which is the same as stroke’s shape in the first two rows, and drawn by users in the last row.

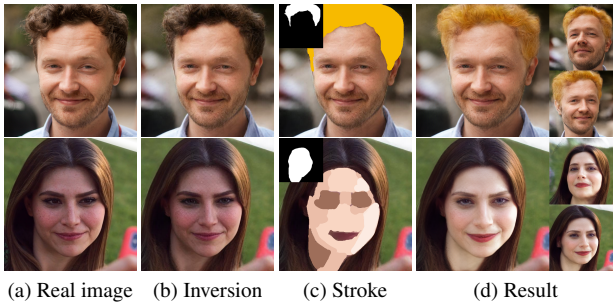


Figure 5. Real image appearance editing results.

R represents the pre-trained facial recognition network, ArcFace [6], and $\langle \cdot, \cdot \rangle$ denotes cosine similarity. Regularization constraints are added to ensure that the edited latent variables are as close as possible to the mean latent variables \bar{w} , making the generated faces more natural.

$$L_{reg} = \|w_{edit} - \bar{w}\|_2 \quad (8)$$

The overall loss function is given by:

$$L = \alpha_1 L_{density} + \alpha_2 L_{img} + \alpha_3 L_{LPIPS} + \alpha_4 L_{ID} + \alpha_5 L_{reg} \quad (9)$$

Here, α_1 , α_2 , α_3 , α_4 and α_5 represent the weighting coefficients for each loss term, with values of 0.05, 2.0, 2.0, 0.2, and 0.005, respectively. In the first stage training, we set $\alpha_1 = 0$ because only appearance is expected to be encoded in this stage.

4. Experiments

Dataset: Our training dataset comprises original faces generated by EG3D random sampling and real images in-

version [1] of CelebA-HQ [20], amounting to 47,134 examples. Appearance-edited faces and corresponding color stroke images are constructed using the method detailed in Section 3.1, resulting in 94,268 paired samples in total.

Training details: We use a single NVIDIA GeForce RTX 3090 GPU to train our framework. The batch size is 4, with the Ranger optimizer [37] and a learning rate of 1e-4. Each training stage involved 30k iterations, completing the final network training process.

Evaluation metrics: We use the following common metrics: Fréchet Inception Distance (FID) [12], Kernel Inception Distance (KID) [4] and Identity Consistency (ID) from [15]. To evaluate the geometry consistency in both 2D and 3D methods, we computed the Structural Similarity Index (SSIM) [35] metric for the edges obtained using the Sobel [8] operator in the edited regions of the generated images compared to the real images, denoted as E-SSIM. A higher E-SSIM value in the edited regions indicates a better preservation of the geometry in those areas.

4.1. Results

Figure 4 illustrates the 3D facial NeRF appearance editing based on color strokes. Benefit from the 3D representation in our framework, users can select arbitrary view point to render images. Then, users could draw color strokes, which are overlaid on extracted strokes extracted based on methods described in Section 3.1 to serve as the edited stroke image. Our framework takes an additional 2D mask inputs, which either have the same shape with drawing strokes, or additionally drawn by users. As shown in Figure 4(c)(e), our method allows the precise control over the skin, hair,

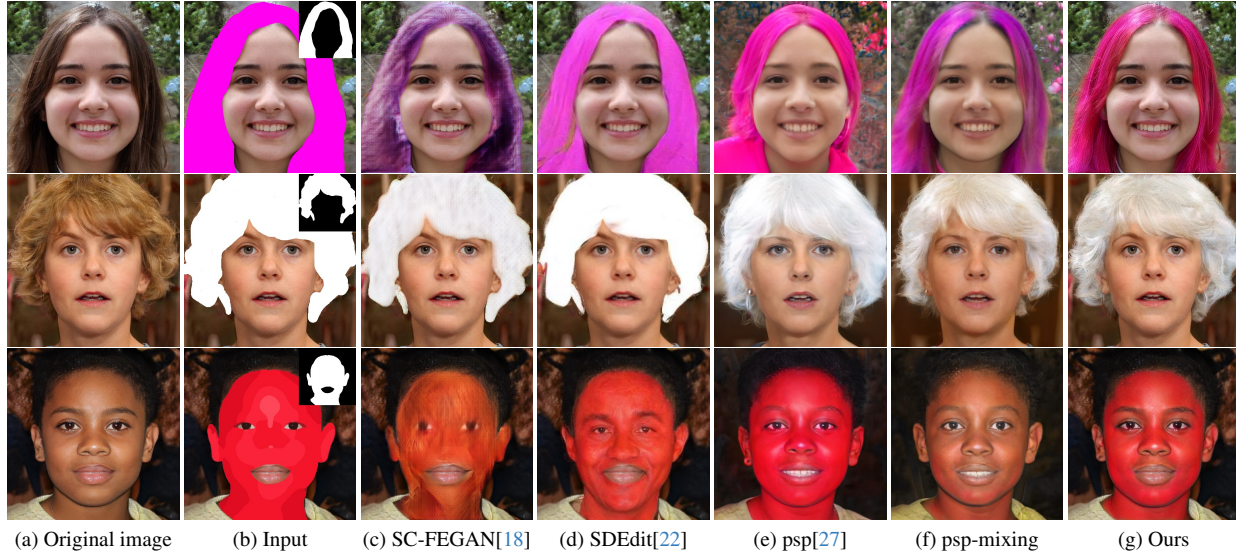


Figure 6. Comparison with existing color stroke-based facial editing methods. SC-FEGAN[18] and SDEdit[22] are 2D image editing approaches that generate texture smoothness. pSp[27] and pSp-mixing globally changes the 3D faces. Our method generates the best results on realism and identity preservation.

and lip colors. It is noteworthy that the geometric details are well preserved as the original faces during appearance editing, as demonstrated by the wrinkle details in the 1st row and hair structure in the 3rd row. Unedited regions, such as hair (1st row) and background, also retain the same as the original faces. The 3rd row of Figure 4 showcases more complex local appearance editing results. To achieve diverse editing effects, users draw finer masks (c) to mark local areas, such as a cluster of hair or specific facial patterns. The 3D local fusion algorithm enables editing in marked areas, resulting in more complex and diverse editing effects.

As shown in Figure 5, given a real facial image, the original facial image is reconstructed in latent space using the PTI algorithm [28] (b). Similarly, users draw color stroke images (c) and edit area masks (d). Our method generates corresponding facial NeRF appearance editing results and obtains multi-view facial outcomes.

4.2. Comparisons

Table 1. Quantitative evaluation using FID, $KID \times 10^2$, mean and std value of identity consistency (ID), $E\text{-SSIM} \times 10^2$ on baseline models and our methods.

	FID (\downarrow)	KID (\downarrow)	ID (\uparrow)	E-SSIM (\uparrow)
SC-FEGAN [18]	47.18	3.41	0.50 ± 0.37	74.36
SDEdit [22]	37.61	2.52	0.41 ± 0.41	74.52
pSp [27]	40.39	2.49	0.66 ± 0.11	78.94
pSp-mixing	31.32	2.00	0.73 ± 0.09	81.50
Ours	16.86	0.83	0.85 ± 0.10	84.30

Since there are no prior works on 3D-aware stroke-based image synthesis, we retrained the pSp [27] model on EG3D using our dataset. To ensure fair comparisons, we further

applied style-mixing on pSp’s predicted latent code by replacing 7-14 layers with original latent codes. We also compare with 2D image editing methods, including SC-FEGAN [18] and SDEdit [22]. Officially released codebases were used for all comparisons. As shown in Figure 6, given the original faces (a), color strokes and masks (b), since the editing areas are too large to complete, SC-FEGAN [18] generates less realistic results (c), with facial details appearing blurry in larger editing areas. Based on diffusion prior knowledge, SDEdit [22] generates more realistic results (d). However, the texture is still too smooth due to the geometry information absence of color strokes. Different from 2D image completion frameworks, pSp [27] predicts a latent code, which, however, totally modify the geometry compared with original faces. Unedited regions, such as background, are also changed. Even utilizing the style-mixing approach (f), there are also undesirable changes and the degree of editing has been largely affected. Compared with existing approaches, our method generates more realistic faces, which shows good faithfulness to the color strokes and precisely preserve the geometry details, even the sophisticated curly hair structure is well maintained.

Quantitative Comparison. To further validate the superiority of our method, we conduct a quantitation evaluation. We generated 1000 facial editing examples in FFHQ datasets, following the data generation method outlined in Section 3.1. We calculated the relevant metrics between all algorithm-generated results and real images. As shown in Table 1, our method has achieved the best results on all metrics, which proves that our method generates more realistic faces (FID, KID). The original identity feature (ID) and geometry details (E-SSIM) are also better preserved.

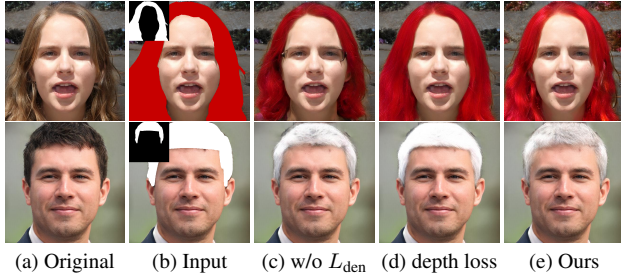


Figure 7. The results of ablation study for density loss. *w/o* L_{den} method is trained without geometry density loss while depth loss method is trained with an alternative depth image loss.

4.3. Ablation Study

Table 2. Ablation study comparison. All models are trained on the datasets described in Section 3.1 and evaluated on 1000 randomly sampled EG3D data. Metric calculations follow the methods outlined in Table 1.

	FID (↓)	KID(↓)	ID(↑)	E-SSIM(↑)
w/o density loss	21.44	1.07	0.67 ± 0.22	77.35
w/ depth loss	23.14	1.23	0.75 ± 0.17	80.10
w/o mask fusion	28.33	1.72	0.82 ± 0.07	82.97
Ours	16.86	0.83	0.85 ± 0.10	84.30

We conduct ablation experiments on EG3D random sampled data to justify the necessity of each component. To validate the effectiveness of the proposed geometry density loss, we presents the results of training without density loss L_{den} and with an alternative depth image loss. Specifically, we employed MSE Loss between the original depth image and the generated one to constrain the network. The results in Figure 7 (d)(e) and Table 2 demonstrate the superiority of our density loss in preserving facial geometry features and maintaining identity consistency.

Figure 8 indicates our volumetric local fusion module’s effectiveness in maintaining consistency for non-edited areas between and the original and edited face. As shown in Figure 8 (d), without using 3D local fusion, non-edited areas undergo significant changes, such as brightening hair in the first row and reddening the face in the second row. In contrast, our method (Figure 8 (e)) perfectly maintains consistency in non-edited areas without affecting the editing effects in edited regions. Figures 8 (d) and (f) display computed heatmap differences between pre-edited and post-edited images.

5. Conclusion

In this study, we propose a facial NeRF appearance editing method based on color strokes, enabling the modification of facial appearance features while preserving the original facial geometry. Initially, a data construction approach is introduced, leveraging StyleGAN2 and semantic segmentation networks to create paired data with consistent

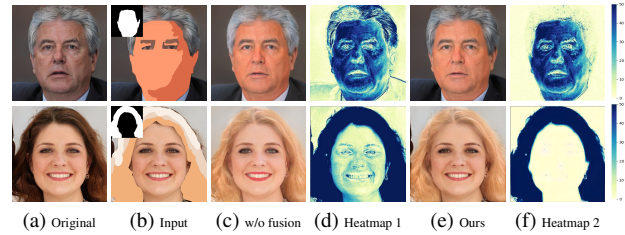


Figure 8. The results of ablation study for local region fusion. *w/o* fusion are the results without 3D local fusion method. Corresponding heatmap indicates the difference between the edited images and the original images.

geometry and diverse appearances, along with corresponding color strokes. Subsequently, an appearance editing network is constructed within the pre-trained generative network’s latent space. The network employs a color stroke encoder to extract color features from input stroke images and utilizes a latent variable fusion network to achieve appearance editing while preserving geometric features. To maintain unchanged regions, a 3D local fusion method is introduced, enabling effective editing of localized appearance while preserving the facial identity features. Experimental results demonstrate the superiority of our method in generating higher-quality digital facial results compared to existing approaches.

Despite achieving high-quality color editing results through color strokes, our method faces several challenges. Firstly, for highly detailed areas such as eyeballs, our method struggles to yield effective editing results. Handling these regions separately through segmentation and dedicated data processing could enhance editing outcomes. Additionally, due to limitations in the training dataset, our method encounters difficulties in handling overly complex color strokes, such as rainbow-colored hair, tending to average color strokes instead. To address this, generating distinct colors separately, followed by mask creation and blending, could yield superior results. Future research endeavors could extend the proposed algorithm from 3D digital faces to other types of 3D models, achieving more generalized editing effects. Moreover, combining color strokes with line drawings and semantic segmentation maps could enable simultaneous editing of 3D digital facial geometry and appearance.

Acknowledgement

This work was supported by National Natural Science Foundation of China (No. 62102403), and the Open Project Program of State Key Laboratory of Virtual Reality Technology and Systems, Beihang University (No. VR-LAB2022C07). We would like to thank Professor Lin Gao for his helpful discussions and generous support.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4432–4441, 2019. 6
- [2] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Restyle: A residual-based stylegan encoder via iterative refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6711–6720, 2021. 4
- [3] David Bau, Hendrik Strobelt, William S. Peebles, Jonas Wulff, Bolei Zhou, Jun-Yan Zhu, and Antonio Torralba. Semantic photo manipulation with a generative image prior. *ACM Transactions on Graphics*, 38(4):59:1–59:11, 2019. 2
- [4] Mikolaj Binkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD gans. In *6th International Conference on Learning Representations*, 2018. 6
- [5] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16102–16112, 2022. 1, 2, 3
- [6] Jiankang Deng, Jia Guo, Xue Niannan, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4690–4699, 2019. 6
- [7] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10663–10673, 2022. 3
- [8] Richard O Duda, Peter E Hart, et al. *Pattern classification and scene analysis*. Wiley New York, 1973. 6
- [9] Lin Gao, Feng-Lin Liu, Shu-Yu Chen, Kaiwen Jiang, Chun-Peng Li, Yu-Kun Lai, and Hongbo Fu. Sketchfacenerf: Sketch-based facial generation and editing in neural radiance fields. *ACM Transactions on Graphics*, 42(4):159:1–159:17, 2023. 3
- [10] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. 1
- [11] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d aware generator for high-resolution image synthesis. In *International Conference on Learning Representations*, 2022. 1
- [12] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6
- [13] Xueqi Hu, Qiusheng Huang, Zhengyi Shi, Siyuan Li, Changxin Gao, Li Sun, and Qingli Li. Style transformer for image inversion and editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11337–11346, 2022. 4
- [14] Xun Huang and Serge J. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *International Conference on Computer Vision*, pages 1510–1519, 2017. 3
- [15] Yuge Huang, Yuhan Wang, Ying Tai, Xiaoming Liu, Pengcheng Shen, Shaoxin Li, and Feiyue Huang Jilin Li. Curricularface: Adaptive curriculum learning loss for deep face recognition. pages 5900–5909, 2020. 6
- [16] Kaiwen Jiang, Shu-Yu Chen, Feng-Lin Liu, Hongbo Fu, and Lin Gao. Nerffaceediting: Disentangled face editing in neural radiance fields. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. 2, 3
- [17] Kyungmin Jo, Gyumin Shim, Sanghun Jung, Soyoung Yang, and Jaegul Choo. Cg-nerf: Conditional generative neural radiance fields for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 724–733, 2023. 2, 3
- [18] Youngjoo Jo and Jongyoul Park. Sc-fegan: Face editing generative adversarial network with user’s sketch and color. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1745–1753, 2019. 2, 7
- [19] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, 2016. 5
- [20] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *6th International Conference on Learning Representations*, 2018. 6
- [21] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8107–8116, 2020. 2, 3
- [22] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *The Tenth International Conference on Learning Representations*, 2022. 2, 7
- [23] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1
- [24] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. 1, 3
- [25] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. Stylesdf: High-resolution 3d-consistent image and geometry generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13503–13513, 2022. 3
- [26] Tiziano Portenier, Qiyang Hu, Attila Szabó, Siavash Arjomand Bigdeli, Paolo Favaro, and Matthias Zwicker.

- Faceshop: deep sketch-based face image editing. *ACM Transactions on Graphics*, 37(4):99, 2018. 2
- [27] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. Encoding in style: a stylegan encoder for image-to-image translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2287–2296, 2021. 2, 4, 5, 7
- [28] Daniel Roich, Ron Mokady, Amit H. Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Transactions on Graphics*, 42(1):6:1–6:13, 2023. 4, 7
- [29] Patsorn Sangkloy, Jingwan Lu, Chen Fang, Fisher Yu, and James Hays. Scribbler: Controlling deep image synthesis with sketch and color. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5400–5409, 2017. 2
- [30] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020. 1, 3
- [31] Jingxiang Sun, Xuan Wang, Yichun Shi, Lizhen Wang, Jue Wang, and Yebin Liu. Ide-3d: Interactive disentangled editing for high-resolution 3d-aware portrait synthesis. *ACM Transactions on Graphics*, 41(6):1–10, 2022. 2, 3
- [32] Jingxiang Sun, Xuan Wang, Yong Zhang, Xiaoyu Li, Qi Zhang, Yebin Liu, and Jue Wang. Fenerf: Face editing in neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7672–7682, 2022. 2, 3
- [33] Carlo Tomasi and Roberto Manduchi. Bilateral filtering for gray and color images. In *Sixth international conference on computer vision (IEEE Cat. No. 98CH36271)*, pages 839–846. IEEE, 1998. 4
- [34] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for stylegan image manipulation. *ACM Transactions on Graphics*, 40(4):1–14, 2021. 2
- [35] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 6
- [36] Chu-Feng Xiao, Deng Yu, Xiaoguang Han, Youyi Zheng, and Hongbo Fu. Sketchhairsalon: deep sketch-based hair image synthesis. *ACM Transactions on Graphics*, 40(6):216:1–216:16, 2021. 2
- [37] Hongwei Yong, Jianqiang Huang, Xiansheng Hua, and Lei Zhang. Gradient centralization: A new optimization technique for deep neural networks. In *Computer Vision European Conference*, pages 635–652. Springer, 2020. 6
- [38] Changqian Yu, Changxin Gao, Jingbo Wang, Gang Yu, Chunhua Shen, and Nong Sang. Bisenet v2: Bilateral network with guided aggregation for real-time semantic segmentation. *International Journal of Computer Vision*, 129:3051–3068, 2021. 2, 4