

DESIGNING A 3D-AWARE STYLENeRF ENCODER FOR FACE EDITING

Songlin Yang^{1,2}, Wei Wang^{2,*}, Bo Peng², Jing Dong²

¹School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China

²Center for Research on Intelligent Perception and Computing, NLPR, CASIA, Beijing, China

yangsonglin2021@ia.ac.cn, {wwang, bo.peng, jdong}@nlpr.ia.ac.cn

ABSTRACT

GAN inversion has been exploited in many face manipulation tasks, but 2D GANs often fail to generate multi-view 3D consistent images. The encoders designed for 2D GANs are not able to provide sufficient 3D information for the inversion and editing. Therefore, 3D-aware GAN inversion is proposed to increase the 3D editing capability of GANs. However, the 3D-aware GAN inversion remains under-explored. To tackle this problem, we propose a 3D-aware (**3Da**) encoder for GAN inversion and face editing based on the powerful StyleNeRF model. Our proposed **3Da** encoder combines a parametric 3D face model with a learnable detail representation model to generate geometry, texture and view direction codes. For more flexible face manipulation, we then design a dual-branch StyleFlow module to transfer the StyleNeRF codes with disentangled geometry and texture flows. Extensive experiments demonstrate that we realize 3D consistent face manipulation in both facial attribute editing and texture transfer. Furthermore, for video editing, we make the sequence of frame codes share a common canonical manifold, which improves the temporal consistency of the edited attributes.

Index Terms— Neural Radiance Field (NeRF), GAN Inversion, 3D Consistent Face Manipulation

1. INTRODUCTION

Face editing via GAN (Generative Adversarial Network) inversion [1] enables users to flexibly edit a wide range of facial attributes in real face images. Existing methods [2, 3, 4, 5] first invert face images into the latent space of 2D GANs such as StyleGAN [6], then manipulate the style codes, and finally feed the edited codes into the pre-trained generator to obtain the edited face images. However, 2D GANs lack the knowledge of the underlying 3D structure of the faces, and their 3D consistency in multi-view generation is limited, as shown in Fig. 1.

In order to increase the 3D consistency of the generators in the GAN-inversion-based manipulation pipeline, one intuitive idea is replacing the 2D GANs with 3D-aware GANs [11, 12, 13, 14, 15]. However, the vanilla encoders designed for 2D GAN inversion fail to provide sufficient 3D information for the 3D-aware GAN inversion. Furthermore, the SOTA 2D encoders like e4e [16] bring much variety in the inversion stage, which degrades the video consistency in video editing. Therefore, to obtain better 3D consistency in multi-view facial attribute editing, we propose a 3D-aware (**3Da**) StyleNeRF [13] encoder which encodes geometry and texture separately to have more flexible manipulation capability.

Our proposed **3Da** encoder combines a parametric 3D face model with a learnable detail representation model to generate the

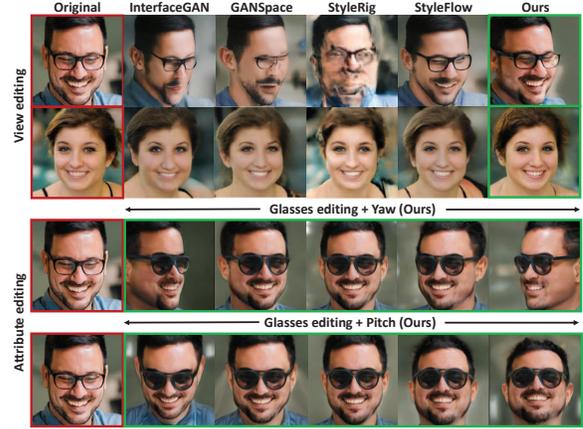


Fig. 1. The comparisons of the multi-view editing effects, using StyleRig [7], InterfaceGAN [8], GANSpace [9], StyleFlow [10] and ours. Our method not only achieves better results in novel views, but also preserves the multi-view 3D consistency of the edited results.

geometry, texture and view direction codes. By introducing the 3D face model, we can enhance the stability of the generated faces. The 3D-aware inversion codes are then fed into a well trained dual-branch StyleFlow [10] module which makes for the flexible face manipulation. We realize the 3D consistent face manipulation in both facial attribute editing and texture transfer. Moreover, we extend our pipeline to video editing. We make the video frames share a common canonical representation manifold, which improves the temporal consistency of the edited attributes.

The main contributions of this work are as follows: We propose the first 3D-aware (**3Da**) StyleNeRF encoder for the face editing. Our **3Da** encoder is able to encode geometry, texture and view direction information separately, achieving multi-view generation and facial attribute editing simultaneously. By introducing the parametric 3D face model, we are able to enhance the stability of the generated faces, which aligns the facial details with the morphable model adaptively.

2. METHOD

2.1. 3D-Aware StyleNeRF Inversion for Face Embedding

The StyleNeRF [13] is adopted as our pre-trained generator. It has two inputs for conditioning the style and camera view respectively. Its NeRF-based [17] architecture performs volume rendering only to produce a low-resolution feature map, and progressively applies up-

* Corresponding author.

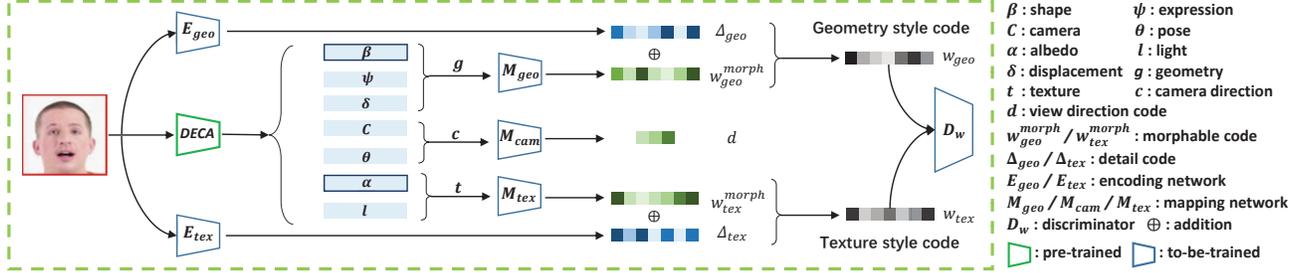


Fig. 2. The framework of our 3D-aware (3Da) encoder. Note that when embedding the video frames, the shape β and albedo α should be the same for the same face among different frames, i.e., frame-irrelevant. Therefore, in the video setting, we first extract these two coefficients for all the frames. Then we use the averaged β and α for encoding all the frames.

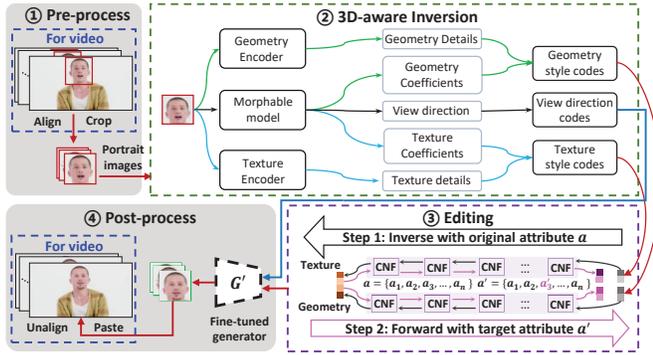


Fig. 3. The portrait manipulation pipeline with our 3Da StyleNeRF encoder for the image and video setting. Our 3D-aware GAN inversion module realizes the disentanglement of canonical and morphable modulations, as well as the separate editing of geometry and texture. Our attribute editing module extends the StyleFlow [10] to two branches for texture and geometry, respectively. Note that our 3Da encoder is able to disentangle the frame-irrelevant information when encoding the video frame code sequence, which is beneficial to increase temporal consistency of the video editing.

sampling in 2D to obtain high-resolution images. Our **first** motivation is that, two branches of geometry and texture should be adopted to match the 3D-aware architecture. The **second** motivation is the disentanglement of canonical and morphable information. So we adopt the parametric 3D face model DECA [18] as 3D prior and use the ResNet-based [19] encoder to provide the detail information. This has two benefits: (1) For training, the sampled style codes with their corresponding synthetic images are used to train the encoder for encoding the images into the GAN space. Our proposed methods can accelerate the convergence and avoid overfitting to the synthetic data. (2) For the video setting, this can guide every code in the code sequence of video frames to share a canonical representation manifold of the target face, preserving temporal consistency of the inversion and editing.

Encoder. As shown in Fig. 2, for the geometry style code, $w_{geo} = w_{geo}^{morph} + \Delta_{geo}$. For the texture style code, $w_{tex} = w_{tex}^{morph} + \Delta_{tex}$. The addition operation combines the 3D information and content details. The CNN-based encoding networks E_{geo} and E_{tex} are used to extract the detail codes Δ_{geo} and Δ_{tex} respectively. For the morphable codes w_{geo}^{morph} , w_{tex}^{morph} and view direction code d , DECA [18] is used to extract the semantic feature vectors

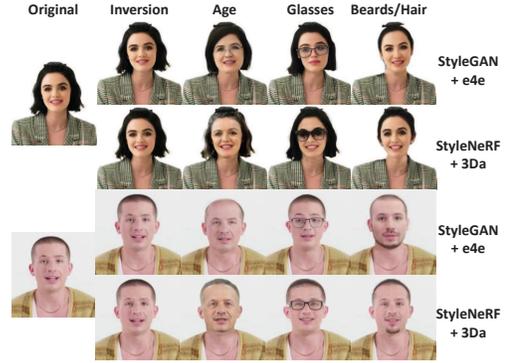


Fig. 4. The edited face images of different attributes. Zoom in the digital version for better view.

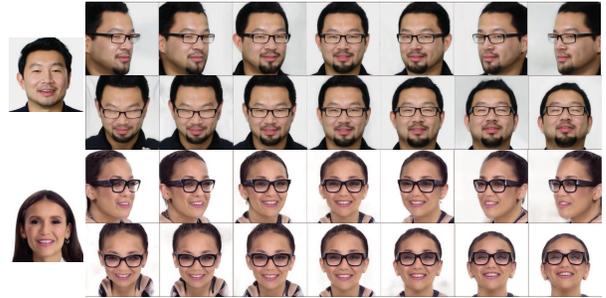


Fig. 5. The multi-view generation of multi-attribute editing.

of geometry g , texture t and camera direction c . Then, geometry g and texture t are input into fully-connected mapping networks M_{geo} and M_{tex} to obtain the morphable codes of geometry w_{geo}^{morph} and texture w_{tex}^{morph} . Camera direction c is input into M_{cam} to get the view direction codes d . Specifically, geometry g is concatenated by shape β and expression ψ and displacement δ . Texture t is concatenated by albedo α , light l . The camera direction c is concatenated by camera C and pose θ .

Discriminator. To encourage the style codes to lie within the distribution of the latent style code space of StyleNeRF, denoted as \mathcal{W} , a discriminator D_w is used to discriminate between real samples from the \mathcal{W} space and the learned latent space of our 3Da encoder. This discriminator is important because it is able to not only accelerate convergence, but also avoid the mode collapse (see Fig. 9).

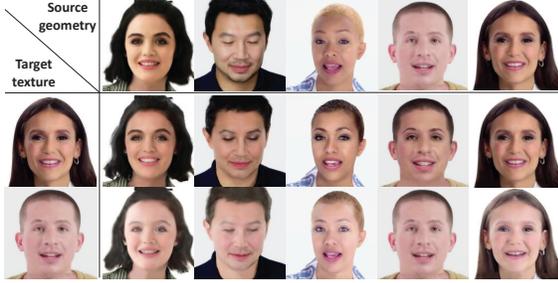


Fig. 6. Texture transfer. The images in the first row provide the source geometry, and the images in the first column provide the target texture. Our **3Da** encoder is able to transfer the target texture to the source geometry.

Method	Age \uparrow	Glasses \uparrow	Beards \uparrow	Hair \uparrow
StyleGAN + e4e	0.637	0.653	0.651	0.694
StyleNeRF + 3Da	0.794	0.791	0.803	0.811

Table 1. The identity consistency scores of edited face images.

Loss Function. For formulation, we denote $E_w(\mathbf{x}) = \{\mathbf{w}_i\}_{1 \leq i \leq N}$ as style codes and $E_d(\mathbf{x}) = \mathbf{d}$ as view direction codes, where N is the number of style modulation layers ($N = 21$ for StyleNeRF). The E_w and E_d represent the networks of our **3Da** encoder. Note that each code in $\{\mathbf{w}_i\}_{1 \leq i \leq 7}$ is the same and referred to as \mathbf{w}_{geo} , while each code in $\{\mathbf{w}_i\}_{8 \leq i \leq 21}$ is the same and referred to as \mathbf{w}_{tex} . To optimize our encoder and discriminator in an adversarial manner, we use the non-saturating GAN loss function [20] to train these networks as follows:

$$\mathcal{L}_{adv}^{D,E} = - \mathbb{E}_{\mathbf{w} \sim W} [\log D_w(\mathbf{w})] - \mathbb{E}_{\mathbf{x} \sim p_X} [\log(1 - D_w(E_w(\mathbf{x})))] \quad (1)$$

$$\mathcal{L}_{rec}^E = \mathcal{L}_{sim} + \lambda_1 \mathcal{L}_{style} + \lambda_2 \mathcal{L}_{view} \quad (2)$$

where \mathcal{L}_{sim} , \mathcal{L}_{style} and \mathcal{L}_{view} are as follows:

$$\mathcal{L}_{sim} = \|\mathbf{x} - G(E_w(\mathbf{x}), E_d(\mathbf{x}))\|_2 + vgg(\mathbf{x}, G(E_w(\mathbf{x}), E_d(\mathbf{x}))), \quad (3)$$

$$\mathcal{L}_{style} = \|\mathbf{w}_{geo} - \mathbf{w}_{geo}^{GT}\|_1 + \|\mathbf{w}_{tex} - \mathbf{w}_{tex}^{GT}\|_1, \quad (4)$$

$$\mathcal{L}_{view} = \|\mathbf{d} - \mathbf{d}^{GT}\|_1, \quad (5)$$

The target image \mathbf{x} is the style-mixing image with the ground-truth geometry style code \mathbf{w}_{geo}^{GT} , texture style code \mathbf{w}_{tex}^{GT} and view direction code \mathbf{d}^{GT} . The G is the fixed pre-trained generator. The vgg denotes perceptual loss [21]. We set λ_1 and λ_2 as 0.5 and 5.

2.2. Dual-Branch StyleFlow for Face Editing

We adopt StyleFlow [10] as the attribute editing method. However, the original StyleFlow only has a single branch of Continuous Normalizing Flow (CNF) blocks, failing to fully utilize the advantages of our **3Da** encoder. Therefore, as shown in the Fig. 3, we train two branches of Continuous Normalizing Flows $\{\phi_s\}_{s=geo,tex}$. Note that ϕ_{geo} and ϕ_{tex} are used to obtain geometry style code \mathbf{w}_{geo} and texture style code \mathbf{w}_{tex} respectively, for controllable editing. We denote \mathbf{v} as the variable of the given StyleNeRF space, while t is the time variable. We suppose that \mathbf{w}_{geo} and \mathbf{w}_{tex} are mapped from a latent variable \mathbf{z} in a normal distribution. We use $\{\phi_s\}_{s=geo,tex}$ to conduct the inversion inference as follows:

$$\mathbf{v}(t_0) = \mathbf{v}(t_1) + \int_{t_1}^{t_0} \phi_s(\mathbf{v}(t), t, \mathbf{a}) dt, \quad (6)$$

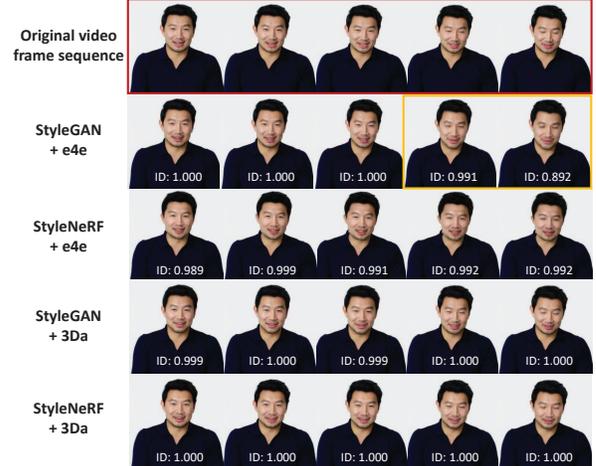


Fig. 7. The inversion of video frame sequences. The StyleNeRF with **3Da** (fifth row) achieves better ID preservation results than others. The StyleGAN with e4e (second row) generates some inversion frames with low ID similarity (marked yellow). The e4e fails to keep the consistency of StyleNeRF, as shown in the third row.

where $\mathbf{v}(t_0)$ is the \mathbf{z} . For $s = geo$, the $\mathbf{v}(t_1)$ is \mathbf{w}_{geo} , while $\mathbf{v}(t_1)$ is \mathbf{w}_{tex} if $s = tex$. Note that \mathbf{a} is the original attribute vector. Then, we modify \mathbf{a} according to the given editing instruction, to obtain the edited attribute vector \mathbf{a}' . After that, we perform a forward inference to produce the edited style code $\mathbf{w}'_{geo} = \mathbf{v}(t_1)$ or $\mathbf{w}'_{tex} = \mathbf{v}(t_1)$, conditioned on \mathbf{a}' as follows:

$$\mathbf{v}(t_1) = \mathbf{v}(t_0) + \int_{t_0}^{t_1} \phi_s(\mathbf{v}(t), t, \mathbf{a}') dt, \quad (7)$$

The above is the inference process, and the training details of CNF blocks can be found in this work [10].

3. EXPERIMENTS

3.1. Implementation Details

Network Architectures. ResNet [19] is used as the backbone for encoding networks E_{geo} and E_{tex} , to extract the feature vectors, corresponding to the input dimensions of StyleNeRF [13]. M_{geo} and M_{tex} are fully-connected networks with 5 layers, while M_{cam} has 3 layers. The LeakyReLU is selected as the activation function. We conduct all the experiments on one NVIDIA RTX 3090. We conducted some preliminary prototyping using the MindSpore framework during our implementation. Our encoder requires 4 days, while the editing module requires 2 days.

Training Data and Annotation. We randomly sample and save 10,000 groups of style-mixing style codes \mathbf{w} , view direction codes \mathbf{d} and their corresponding StyleNeRF-generated images as training data. Moreover, the attribute vectors of these generated images are annotated using Microsoft Face API [22], which every dimension of the vector represents an attribute. These attribute vectors and their corresponding style codes are used for training our dual-branch StyleFlow-based attribute editing module.

Baseline and Compared Methods. The baseline of our experiments is the StyleGAN [6] with e4e [16] encoder that is widely used in this field. Our differences are as follows: **(1)** Inputs: StyleNeRF has input of style codes and view direction codes, while StyleGAN has

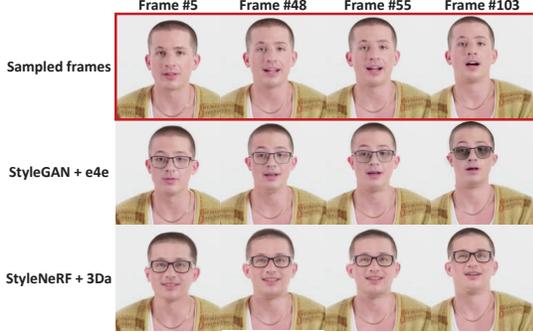


Fig. 8. Comparison of the temporal consistency of video editing. The ‘StyleGAN + e4e’ method generates same changing degree of *Glasses* leads to *Black rimmed glasses* in the previous frames, while *Sun glasses* in the later frames. Our **3Da** with StyleNeRF is able to maintain the temporal consistency of the edited attributes.

Method	PSNR \uparrow	SSIM \uparrow	VIF \uparrow	FVD \downarrow
StyleGAN + e4e	40.3	0.996	0.75	64.2
StyleNeRF + e4e	32.1	0.995	0.71	168.1
StyleGAN + 3Da	38.4	0.995	0.78	121.7
StyleNeRF + 3Da	41.2	0.997	0.79	51.3

Table 2. The quantitative evaluation of generation quality.

only style codes. **(2) Style codes:** The **3Da** encoder has only two different style codes, i.e., geometry and texture codes, while e4e output N different style codes ($N = 21$ or $N = 18$ for StyleNeRF or StyleGAN). **(3) Basic architecture:** StyleNeRF adopts NeRF [17] as basic generation networks, while StyleGAN lacks 3D prior. All the experiments are done under the same setting. We also compare with some other SOTA methods in the experiments, i.e., StyleRig [7], InterfaceGAN [8], GANSpace [9], StyleFlow [10].

Metrics. Quality: PSNR, SSIM and VIF are used to measure the generation quality. Frechet Video Distance (FVD) [23] extends the FID [24] to video quality settings. **Identity Consistency:** We evaluate the identity of the edited faces using ArcFace [25] cosine similarity score. **Attribute Consistency:** Following StyleFlow [10], we use ResNet-18 [19] trained on CelebA [26] as the facial attribute prediction model to output the attribute vectors of face images.

3.2. Face Image Inversion and Editing

Attribute Editing. As shown in Fig. 4, we select *Age*, *Glasses*, *Beards* and *Hair* as the examples. As shown in Tab. 1, we evaluate the identity consistency scores of edited face images compared with their original images in the FFHQ [27] dataset. The generation quality is quantitatively evaluated in Tab. 2.

Multi-Attribute Editing and Multi-View Generation. As shown in Fig. 1, our method has good 3D consistency among different views of edited images. Furthermore, as shown in Fig. 5, our method can handle the multi-view generation and simultaneously edit multiple attributes.

Texture Transfer. As shown in Fig. 6, we can realize texture transfer among different real images by combing the geometry style code of one image with the texture style code of another and then inputting the style-mixing codes to the StyleNeRF. This illustrates the good geometry-texture disentanglement ability of our method.



Fig. 9. Ablation study of inversion results. **(a)** shows **3Da** encoder with discriminator and without real training data. **(b)** is mapping DECA coefficients to style codes. **(c)** is **3Da** encoder without the discriminator. **(d)** is **3Da** encoder with real data and the discriminator. Note that they are trained with the same epochs.

Method	Age \downarrow	Glasses \downarrow	Beards \downarrow	Hair \downarrow
StyleGAN + e4e	0.498	0.533	0.502	0.497
StyleNeRF + 3Da	0.454	0.223	0.385	0.441

Table 3. Temporal attribute inconsistency scores of video editing.

3.3. Portrait Video Manipulation

As shown in Fig. 3, our video manipulation pipeline is composed of three main stages, inspired by the STIT [5] method. First, we use DECA [18] to extract the frame-irrelevant information (shape β and albedo α), encode the cropped face images and smooth the style code sequence over a window of two frames by weighted sum rules. Then, the cropped images and style code sequence are used to fine-tune the StyleNeRF generator. Note that the style code sequence is fixed in this fine-tuning process. And lastly, the style code sequence is input to our dual-branch StyleFlow to obtain an edited style code sequence conditioned on the required attribute vector, and the fine-tuned generator is used to obtain the edited frame sequence.

Video Inversion. As shown in Fig. 7, we evaluate different methods under the same setting as STIT [5]. The StyleNeRF with our **3Da** encoder can achieve better results, in aspects of reconstruction and temporal consistency of identity similarity.

Video Editing. As shown in Fig. 8, our method has more consistent video editing effects. As shown in Tab. 3, we quantitatively measure the Mean Absolute Error of the attribute vectors between the first frame and the following frames on 50 testing videos. Each of them has 100 frames with different attribute edited. Our **3Da** encoder embeds the frame sequence more stably, as shown by the lower temporal attribute inconsistency.

3.4. Ablation Study

Fig. 9 (b) shows that only using DECA coefficients is insufficient. As shown in Fig. 9 (c) and Fig. 9 (d), **3Da** without the discriminator leads to the mode collapse, and using real data in training degrades the image quality. Using the real images without ground-truth style codes for training only has the \mathcal{L}_{sim} loss, which is less effective than the embedding supervision from \mathcal{L}_{style} and \mathcal{L}_{view} .

4. CONCLUSION

In this paper, we propose a 3D-aware (**3Da**) StyleNeRF encoder to encode geometry, texture and view direction of the real face images. Extensive experiments qualitatively and quantitatively demonstrate that, we are able to realize high-quality multi-view generation and facial attribute editing. Moreover, we extend our method to the portrait video manipulation, achieving better temporal consistency over the 2D-GAN-based editing methods.

5. ACKNOWLEDGMENTS

This work is supported by the National Key Research and Development Program of China under Grant No. 2021YFC3320103, the National Natural Science Foundation of China (NSFC) under Grant 61972395, 62272460, a grant from Young Elite Scientists Sponsorship Program by CAST (YESS), and sponsored by CAAI-Huawei MindSpore Open Fund.

6. REFERENCES

- [1] Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, and Ming-Hsuan Yang, “Gan inversion: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [2] Yuming Jiang, Ziqi Huang, Xingang Pan, Chen Change Loy, and Ziwei Liu, “Talk-to-edit: Fine-grained facial editing via dialog,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 13799–13808.
- [3] Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujiu Yang, “Styleheat: One-shot high-resolution editable talking face generation via pretrained stylegan,” *arXiv preprint arXiv:2203.04036*, 2022.
- [4] Rameen Abdal, Peihao Zhu, Niloy J Mitra, and Peter Wonka, “Video2stylegan: Disentangling local and global variations in a video,” *arXiv preprint arXiv:2205.13996*, 2022.
- [5] Rotem Tzaban, Ron Mokady, Rinon Gal, Amit H Bermano, and Daniel Cohen-Or, “Stitch it in time: Gan-based facial editing of real videos,” *arXiv preprint arXiv:2201.08361*, 2022.
- [6] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila, “Analyzing and improving the image quality of stylegan,” in *CVPR*, 2020, pp. 8110–8119.
- [7] Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhofer, and Christian Theobalt, “Stylerig: Rigging stylegan for 3d control over portrait images,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6142–6151.
- [8] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou, “Interfacegan: Interpreting the disentangled face representation learned by gans,” *IEEE transactions on pattern analysis and machine intelligence*, 2020.
- [9] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris, “Ganspace: Discovering interpretable gan controls,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 9841–9850, 2020.
- [10] Rameen Abdal, Peihao Zhu, Niloy J Mitra, and Peter Wonka, “Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows,” *ACM Transactions on Graphics (TOG)*, vol. 40, no. 3, pp. 1–21, 2021.
- [11] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger, “Graf: Generative radiance fields for 3d-aware image synthesis,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 20154–20166, 2020.
- [12] Michael Niemeyer and Andreas Geiger, “Giraffe: Representing scenes as compositional generative neural feature fields,” in *CVPR*, 2021, pp. 11453–11464.
- [13] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt, “Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis,” *ICLR*, 2022.
- [14] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al., “Efficient geometry-aware 3d generative adversarial networks,” in *CVPR*, 2022, pp. 16123–16133.
- [15] Yinghao Xu, Sida Peng, Ceyuan Yang, Yujun Shen, and Bolei Zhou, “3d-aware image synthesis via learning structural and textual representations,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18430–18439.
- [16] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or, “Designing an encoder for stylegan image manipulation,” *ACM Transactions on Graphics (TOG)*, vol. 40, no. 4, pp. 1–14, 2021.
- [17] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” in *ECCV*. Springer, 2020, pp. 405–421.
- [18] Yao Feng, Haiwen Feng, Michael J Black, and Timo Bolkart, “Learning an animatable detailed 3d face model from in-the-wild images,” *ACM Transactions on Graphics (TOG)*, vol. 40, no. 4, pp. 1–13, 2021.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [20] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, “Generative adversarial nets,” *Advances in neural information processing systems*, vol. 27, 2014.
- [21] Justin Johnson, Alexandre Alahi, and Li Fei-Fei, “Perceptual losses for real-time style transfer and super-resolution,” in *European conference on computer vision*. Springer, 2016, pp. 694–711.
- [22] Azure, “<https://azure.microsoft.com/en-in/services/cognitive-services/face/>,” *Microsoft*, 2020.
- [23] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly, “Towards accurate generative models of video: A new metric & challenges,” *arXiv preprint arXiv:1812.01717*, 2018.
- [24] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.
- [25] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou, “Arcface: Additive angular margin loss for deep face recognition,” in *CVPR*, 2019, pp. 4690–4699.
- [26] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang, “Large-scale celebfaces attributes (celeba) dataset,” *Retrieved August*, vol. 15, no. 2018, pp. 11, 2018.
- [27] Tero Karras, Samuli Laine, and Timo Aila, “A style-based generator architecture for generative adversarial networks,” in *CVPR*, 2019, pp. 4401–4410.