

# DEFORMTOON3D: Deformable Neural Radiance Fields for 3D Toonification

Junzhe Zhang<sup>1,3\*</sup> Yushi Lan<sup>1\*</sup> Shuai Yang<sup>1</sup> Fangzhou Hong<sup>1</sup>  
 Quan Wang<sup>3</sup> Chai Kiat Yeo<sup>2</sup> Ziwei Liu<sup>1</sup> Chen Change Loy<sup>1</sup>

<sup>1</sup>S-Lab, Nanyang Technological University

<sup>2</sup>Nanyang Technological University <sup>3</sup>SenseTime Research

{shuai.yang, asckyeo, ziwei.liu, ccloy}@ntu.edu.sg

{junzhe001, yushi001, fangzhou001}@e.ntu.edu.sg {wangquan}@sensetime.com

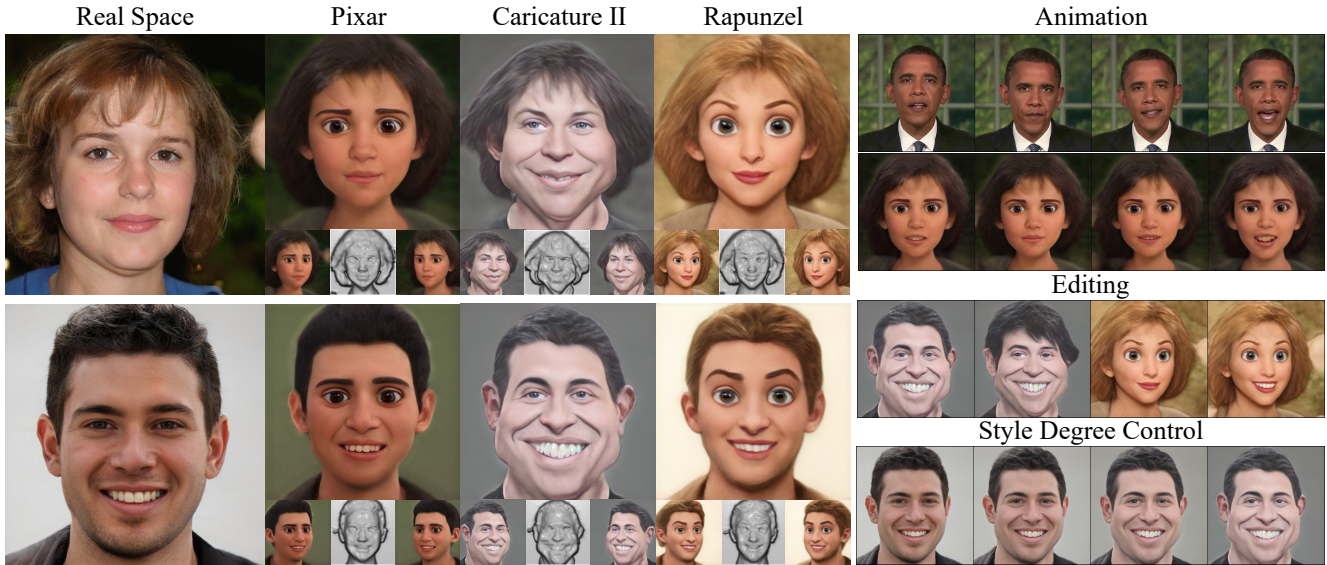


Figure 1: We propose DEFORMTOON3D, an efficient 3D toonification framework which supports geometry-texture decomposed toonification over multiple styles. DEFORMTOON3D is fine-tuning free and could be easily extended to a series of downstream applications designed for the pre-trained GAN, *e.g.*, animation given a driving video, semantic attributes editing (+ bangs and smile), and flexible style degree control.

## Abstract

In this paper, we address the challenging problem of 3D toonification, which involves transferring the style of an artistic domain onto a target 3D face with stylized geometry and texture. Although fine-tuning a pre-trained 3D GAN on the artistic domain can produce reasonable performance, this strategy has limitations in the 3D domain. In particular, fine-tuning can deteriorate the original GAN latent space, which affects subsequent semantic editing, and requires independent optimization and storage for each new style, limiting flexibility and efficient deployment. To overcome these challenges, we propose DEFORMTOON3D, an effective toonification framework tailored for hierarchical 3D GAN. Our approach decomposes 3D toonification into

subproblems of geometry and texture stylization to better preserve the original latent space. Specifically, we devise a novel StyleField that predicts conditional 3D deformation to align a real-space NeRF to the style space for geometry stylization. Thanks to the StyleField formulation, which already handles geometry stylization well, texture stylization can be achieved conveniently via adaptive style mixing that injects information of the artistic domain into the decoder of the pre-trained 3D GAN. Due to the unique design, our method enables flexible style degree control and shape-texture-specific style swap. Furthermore, we achieve efficient training without any real-world 2D-3D training pairs but proxy samples synthesized from off-the-shelf 2D toonification models. Code is released at <https://github.com/junzhezhang/DeformToon3D>.

\*Equal contribution.

## 1. Introduction

Artistic portraits are prevalent in various applications such as comics, animation, virtual reality, and augmented reality. In this work, our main objective is to propose an effective approach for 3D-aware artistic toonification, a critical problem that involves transferring the style of an artistic domain onto a target 3D face with stylized geometry and texture. The task opens up potential applications for quick high-quality 3D avatar creation based on a photograph with the style of a designated artwork, which would typically require highly professional handcraft skills.

Substantial progress has been made in automatic portrait style transfer over 2D images. Starting with image style transfer [15, 55, 35, 29] and image-to-image translation [28, 34, 10, 6, 58], recent advancements in StyleGAN-based generators [24, 25] have shown their potential in high-quality toonification via efficient transfer learning [49]. Specifically, a pre-trained StyleGAN generator on face images is fine-tuned to transfer to the artistic portrait domain. With the progress of 3D-aware GANs [8, 41], researchers have extended this pipeline to 3D with well-designed domain adaptation frameworks [22, 32, 27, 1], enabling remarkable 3D portrait toonification.

Although fine-tuning the pre-trained StyleGAN-based model for toonification achieves superior quality, it has several limitations. **First**, fine-tuning the pre-trained generator shifts its generative space from the real face domain to the artistic portrait domain at the cost of deteriorating the original GAN latent space. With tremendous off-the-shelf tools [56, 64] trained for the original GAN space, altering the well-learned style space would affect the performance of downstream applications over the toonified portrait, *e.g.*, semantic editing. **Second**, despite fine-tuning-based domain adaptations have been thoroughly investigated for 2D GANs [50, 68], applying this technique to 3D GANs fails to leverage the full potential of the architecture of 3D generator models [8, 41] for characterizing view-consistent shape and high-frequency textures in the artistic domain. **Third**, it is inevitable to fine-tune a heavy generator for each new style, which requires hours of training time and additional storage. This limitation affects scalability when deploying dozens of fine-tuned generators for real-time user interactions. Therefore, 3D toonification remains a challenging task that requires further exploration.

To better preserve the pre-trained GAN latent space and to better exploit the 3D GAN generator, we propose **DEFORMTOON3D** that decomposes geometry and texture stylization into more manageable subproblems. In particular, unlike conventional 3D toonification approaches that fine-tune the whole 3D GAN generator following existing 2D fine-tuning schemes, we carefully consider the characteristics of 3D GANs to decompose the stylization of geometry and texture domains. To achieve geometry stylization,

we introduce a novel **StyleField** on top of a pre-trained 3D generator to deform each point in the style space to the pre-trained real space guided by an instance code. This allows for easy extension to multiple styles with a single stylization field by introducing a style code to guide the deformation. Since StyleField already handles geometry stylization well, texture stylization can be easily achieved through adaptive style mixing which injects artistic domain information into the network for effective texture toonification. Notably, our unique design enables training of the method at minimal cost using synthetic paired data with realistic faces generated by a pre-trained 3D GAN and corresponding paired stylized data generated by an off-the-shelf 2D toonification model [70].

The proposed DEFORMTOON3D achieves high-quality geometry and texture toonification over a vast variety of styles, as demonstrated in Fig. 1. Additionally, our approach preserves the original GAN latent space, enabling compatibility with existing tools built on the real face space GAN, including inversion [32], editing [56], and animation [64]. Furthermore, our design significantly reduces the storage footprint by requiring only a small stylization field with a set of AdaIN parameters for artistic domain stylization. In summary, our work makes the following contributions:

- We propose a novel StyleField that separates geometry toonification from texture, providing a more efficient method for modeling 3D shapes than fine-tuning and enabling flexible style control.
- We present an approach to achieve multi-style toonification with a single model, facilitating cross-style manipulation and reducing storage footprint.
- We introduce a full synthetic data-driven training pipeline that offers an efficient and cost-effective solution to training the model without requiring real-world 2D-3D training pairs.

## 2. Related Work

**3D Generative Models.** Inspired by the success of Generative Adversarial Networks (GAN) [16] in generating photo-realistic images [24, 5, 26], researchers have been making efforts towards 3D-aware generation [38, 19, 42]. Starting with explicit intermediate shape representations, such as voxels [38, 19] and meshes [42], which lack photo-realism and are memory-inefficient, researchers have recently shifted towards using implicit functions [44, 36, 9] along with physical rendering processes [60, 37] as intrinsic 3D inductive biases. Among these approaches, 3D generative models [7, 54] extended from neural radiance fields (NeRF) [37] have demonstrated impressive view-consistency in synthesized results. While the original NeRF

is limited to modeling static scenes, recent research has introduced deformation fields to enable NeRF to model dynamic volumes [45, 46, 65, 31]. To increase the resolution of generated images, recent studies [8, 20] have resorted to voxel-based representations or adopted a hybrid design [39, 41, 8, 17]. This hybrid design involves a cascade model  $G = G_1 \circ G_0$  pairing a 3D generator  $G_0$  with a 2D super-resolution decoder  $G_1$ . Both  $G_0$  and  $G_1$  follow the style-based architecture [24, 25] to accept a latent code  $\mathbf{w}$  to control the style of the generated object. By super-resolving the intermediate low-resolution 2D features produced by the  $G_0$  with the  $G_1$ , the hybrid design achieves view-consistent synthesis at high resolution, *e.g.*, 1024<sup>2</sup>.

**Domain Adaptation for StyleGAN Toonification.** Researchers have typically employed domain adaptation in 2D space to achieve toonification with StyleGANs [24, 26]. Typically, a portrait StyleGAN pre-trained on real images [33, 23] is fine-tuned on an artistic domain dataset to generate toonified faces. Building upon this straightforward framework [50], a series of in-depth research has been conducted to further improve style control [70, 69], choices of latent code [61], few-shot training [40], and text-guided adaptation [14].

Pre-trained 3D GANs offer high-quality generation and thus have the potential to facilitate downstream applications such as portrait stylization through 3D GAN inversion [32, 8]. To extend 2D StyleGAN domain adaptation to 3D, CIPS-3D [75] proposed fine-tuning only the super-resolution decoder module for view-consistent texture toonification. However, it is limited to texture toonification since the 3D generator is left unchanged. Dr.3D [22], E3DGE [32], and 3DAvatarGAN [1] fine-tune the entire 3D generator [8, 41] for both geometry and texture toonification, while DATID-3D [27] leverages a Stable Diffusion [53] generated corpus for text-guided domain adaptation. While these methods yield impressive results, they come with limitations, as they require costly fine-tuning and independent model storage for each new style. Furthermore, previous fine-tuning-based toonification methods suffer from limited generality due to the incompatibility with abundant editing techniques developed for the original StyleGAN latent space [56, 57, 47]. In contrast, DEFORMTOON3D fully preserves the original 3D GAN latent space, which makes it intrinsically compatible with the editing methods trained for the original StyleGAN space. It achieves comparable visual quality while being 10 times more storage-efficient than previous methods.

### 3. DEFORMTOON3D

We present the framework of DEFORMTOON3D in Fig. 2. Our approach begins with a typical hybrid 3D-aware design [41, 8] that generates real-domain faces, and reformulates it to a 3D toonification framework. The approach

starts with a cascade model  $G = G_1 \circ G_0$ , which pairs a 3D generator  $G_0$  with a 2D super-resolution decoder  $G_1$ . The generator  $G_0$  captures the underlying geometry with the instance code  $\mathbf{w}$  and camera pose  $\xi$ , and produces an intermediate feature map  $\mathbf{F}$  with volume rendering [37]. Then,  $G_1$  upsamples  $\mathbf{F}$  to obtain a high-resolution image  $\mathbf{I}$  with high-frequency details added. To adapt  $G$  from the real domain to the artistic or cartoon domain, existing methods [22, 1, 75, 32] view  $G_1$  and  $G_0$  as a whole and simply fine-tune the pre-trained  $G$ , failing to take advantage of the decomposed characteristics of the hybrid framework design. By comparison, DEFORMTOON3D fully exploits this cascaded synthesis process by using a novel StyleField module for geometry stylization, which in turn also benefits appearance stylization, allowing it to adopt a simple adaptive style mixing strategy. In Sec. 3.1, we elaborate the proposed StyleField along with the pre-trained  $G_0$  to handle geometry stylization. In Sec. 3.2, we explain how the adaptive style mixing injects the style of the target domain into  $G_1$  to achieve texture stylization. Lastly, we present our training pipeline in Section 3.3.

#### 3.1. Geometry Toonification with StyleField

To train a 3D generator  $\tilde{G}_0$  capable of synthesizing artistic domain geometry, previous methods [22, 27, 32, 1, 62] fine-tune the pre-trained  $G_0$  with a target-domain dataset, which can be computationally expensive and could potentially deteriorate the original GAN latent space. To address this issue, we propose to establish a correspondence between the stylized NeRF  $\mathcal{N}_S$  and the real-space NeRF  $\mathcal{N}_R$ . More specifically, we use a stylization field (the StyleField),  $H_D$ , to bridge the correspondence, such that  $\tilde{G}_0(\mathbf{x}_S) = (G_0 \circ H_D)(\mathbf{x}_S)$ . As shown in Fig. 2, given a stylized NeRF  $\mathcal{N}_S : \mathbb{R}^3 \mapsto \mathbb{R}^4$  of the target style domain, our goal is to estimate a 3D deformation residual,  $H_D : \mathbb{R}^3 \mapsto \mathbb{R}^3$ , which maps  $\mathcal{N}_S$  back to  $\mathcal{N}_R$  via:

$$\mathcal{N}_S \rightarrow \mathcal{N}_R : \mathbf{x}_R = (\mathbf{x}_S + H_D(\mathbf{x}_S)), \forall \mathbf{x}_S \in \mathcal{N}_S, \quad (1)$$

where  $H_D$  represents the residual 3D deformation  $H_D(\mathbf{x}_S) = \Delta \mathbf{x}_S$  in the 3D space of the Stylized NeRF  $\mathcal{N}_S$  and maps each 3D point  $\mathbf{x}_S \in \mathcal{N}_S$  in the stylized space to its corresponding position in the real space  $\mathcal{N}_R$ .

To improve expressiveness, we extend  $H_D$  as a conditional neural deformation field that outputs the offsets under the conditions of style and identity:

$$\mathcal{N}_S \rightarrow \mathcal{N}_R : \mathbf{x}_R = \mathbf{x}_S + H_D(\mathbf{x}_S, \mathbf{w}_S, \mathbf{w}_R) \quad (2)$$

where  $\mathbf{w}_S$  is the style code that specifies the artistic domain,  $\mathbf{w}_R$  is the instance code corresponding to  $\mathcal{N}_R$  that represents the identity of the 3D face in the source domain. Both  $\mathbf{w}_S$  and  $\mathbf{w}_R$  serve as the holistic geometry indicators to guide the deformation.

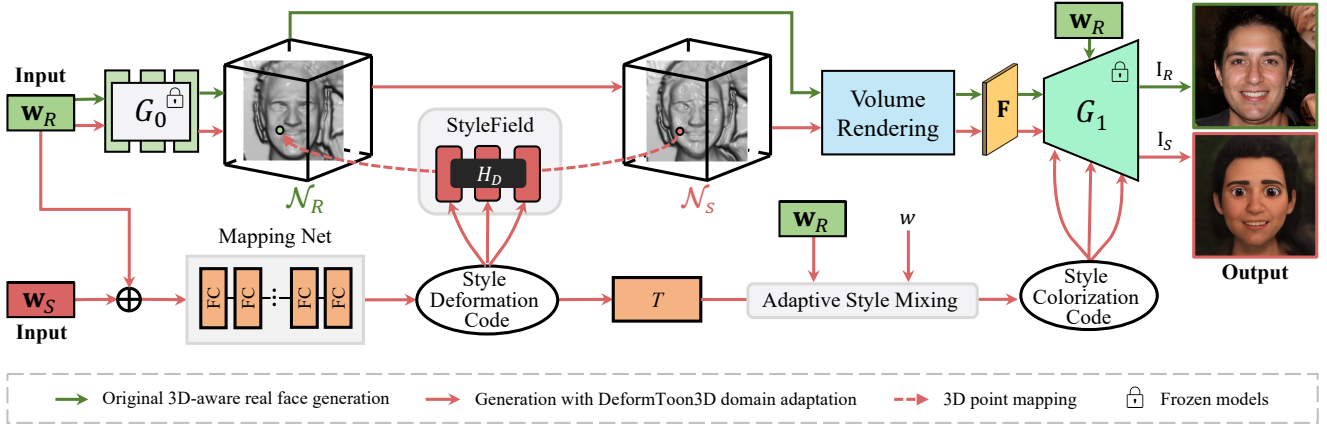


Figure 2: **DEFORMTOON3D framework.** Given sampled instance code  $\mathbf{w}_R$  and style code  $\mathbf{w}_S$  as conditions, DEFORMTOON3D first deforms a point from the style space  $\mathcal{N}_S$  to the real space  $\mathcal{N}_R$ , which achieves geometry toonification without modifying the pre-trained  $G_0$ . Afterwards, we leverage adaptive style mixing with weight  $w$  to inject the texture information of the target domain into the pre-trained  $G_1$  for texture toonification. Both pre-trained generators  $G_0$  and  $G_1$  are kept frozen.

We build  $H_D$  as an MLP consisting of four SIREN [59] layers due to its superior high-frequency modeling capacity. After all the points  $\mathbf{x}_S \in \mathcal{N}_S$  are deformed to the real space, we generate the new feature map  $\hat{\mathbf{F}} = \tilde{G}_0(\mathbf{w}_R, \mathbf{w}_S, \xi)$  for the input of  $G_1$  and synthesize the high-resolution image with the geometry of the artistic domain. By introducing the 3D deformation module  $H_D$ , we no longer need to fine-tune pre-trained  $G_0$  to achieve geometry deformation of the target domain. This greatly alleviates the parameters to optimize by 50% and fully preserves the original GAN latent space. Moreover, with  $\mathbf{w}_S$  serving as the style condition, a single  $H_D$  can support multiple styles, which further saves storage by 98.5% compared to fine-tuning the whole model per style with 10 styles.

### 3.2. Texture Transfer with Adaptive Style Mixing

Thanks to the StyleField  $H_D$  that handles the geometry toonification within the cascade 3D GAN model, we only need to inject the artistic domain texture information into  $G_1$  for texture stylization. Here,  $G_1$  is a 2D style-based architecture, where image styles are effectively adjusted by AdaIN [21]. Inspired by style mixing [24, 26, 70] that controls the AdaIN parameters, we inject the texture information of the target style  $\mathbf{w}_S$  by mixing the style parameters of  $G_1$ , as shown in Fig. 2. To bridge the domain gap between the real space and target domain, we further add a lightweight MLP,  $T$ , for each layer of  $G_1$  to adjust the style code  $\mathbf{w}_S$ . The adapted  $T(\mathbf{w}_S)$  and the instance code  $\mathbf{w}_R$  are fused with a weight  $w$  by weighted average, and sent to the affine transformation block of  $G_1$  to obtain the final style parameters for AdaIN. This mechanism allows us to model and control multi-domain textures with  $T$  and  $w$ , without fine-tuning the original decoder  $G_1$ .

Let  $\tilde{G}_1$  denote  $G_1$  with the adaptive style mixing. The image generation process with domain adaptation given  $w$

and  $\mathbf{w}_S$  can be formulated as  $\hat{\mathbf{I}} = \tilde{G}_1(\hat{\mathbf{F}}, \mathbf{w}_R, \mathbf{w}_S, w)$ .

### 3.3. Training

**Data Preparation.** We follow Sim2Real [67, 74, 32] to generate paired data for training. Specifically, to generate the training corpus for each iteration of the training process, we pre-calculate a set of real space NeRFs  $\mathcal{N}_R$  with corresponding latent codes  $\mathbf{w}_R$  and rendered images  $\mathbf{I}_R$ . To generate pair-wise stylized ground truths, we stylize the rendered image with existing 2D toonification models to obtain the target ground truth  $\mathbf{I}_S$ . Here, to validate our method’s performance on multi-style domain adaptation, we adopt the exemplar-based DualStyleGAN [70] as our 2D toonification model, since it supports hundreds of diverse styles, such as Cartoon, Pixar, and Caricature.

Finally, we define  $\mathcal{X} = \{\mathbf{I}_R, \mathbf{w}_R, \mathbf{w}_S, \mathbf{I}_S\}$  as a training set for DEFORMTOON3D with  $\{\mathbf{w}_R, \mathbf{I}_R, \mathbf{w}_S\}$  to serve as the training inputs and  $\{\mathbf{I}_S\}$  is the set of training ground truth.  $\mathbf{I}_R \in \mathcal{X}$  is drawn i.i.d from distribution  $P(G(\mathbf{z}, \xi))$  where  $\mathbf{z} \sim \mathcal{N}(0, 1)$  and  $\xi$  is the camera pose distribution of pre-trained 3D GAN  $G$ . Some samples are shown in Fig. 3.

**Reconstruction Loss.** Here we use the LPIPS loss [73] to evaluate stylized image quality:

$$\mathcal{L}_{\text{Rec}}(\mathcal{X}) = \mathbb{E}_{\mathcal{X}} \left[ \|P(\hat{\mathbf{I}}) - P(\mathbf{I}_S)\|_2 \right], \quad (3)$$

where  $P(\cdot)$  denotes the perceptual feature extractor.

**Smoothness Regularization.** To encourage the smoothness of stylization field deformation offsets and reduce spatial distortion, a smoothness regularization is included to regularize  $H_D$ . Here we penalize the norm of the Jacobian matrix  $\mathbb{J}_{H_D} = \nabla H_D$  of the deformation field [45] to ensure the learned deformations are physically smooth:

$$\mathcal{L}_{\text{Elastic}} = \text{ReLU}(\|\nabla H_D(\mathbf{x}_S, \mathbf{w}_S, \mathbf{w}_R)\|_2^2 - \epsilon), \quad (4)$$

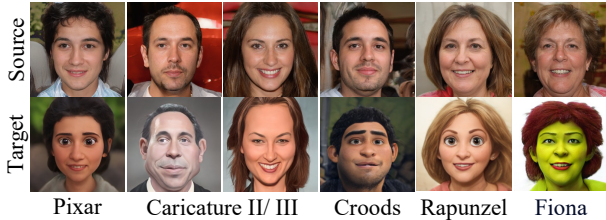


Figure 3: **Training Samples.** We show the source image (row 1) and the stylized images of the target domain (row 2) for training supervisions.

where  $\epsilon$  is the slack parameter for the smoothness regularization. We set  $\epsilon = 0.1$  for all the experiments.

**Adversarial Training.** Additionally, we apply a non-saturating adversarial loss  $\mathcal{L}_{Adv}$  [24] to bridge the domain gap of toonified results.

In summary, the overall loss function is defined as

$$\mathcal{L} = \mathcal{L}_{Rec} + \lambda_{Elastic}\mathcal{L}_{Elastic} + \lambda_{Adv}\mathcal{L}_{Adv},$$

where we set  $\lambda_{Adv} = 0.05$  and  $\lambda_{Elastic} = 0.01$  in all the experiments.

## 4. Experiments

**Datasets.** We mainly focus on the human face domain and use synthesized data for the whole training. We follow the data synthesis procedure mentioned in Sec. 3.3 for generating training corpus and adopt DualStyleGAN [70] as the 2D stylization model to infer paired toonified image. We generate images with the following 10 styles: Pixar, Comic, Slam Dunk, The Croods, Fiona (Shrek), Rapunzel (Disney Princess), Hiccup Horrendous Haddock III (How To Train Your Dragon), and three different caricature styles. We use StyleSDF [41] pre-trained on FFHQ [24] as our 3D generator. For evaluating toonification on real-world images, we evaluate on CelebA-HQ [23].

**Implementation Details.** In all the experiments, we set the learning rate to  $5 \times 10^{-4}$ . We adopt Adam [30] optimizer to train the toonification models. We train the DEFORMTOON3D for 100 epochs. To expedite the training process, we disable the adversarial loss for the initial 50 epochs. The training takes approximately 24 hours using 8 Tesla V100 GPUs with a batch size set to 16. For adaptive style mixing,  $w = 1$  is fixed during training and could be manually selected from  $[0, 1]$  to achieve texture interpolation during inference. More implementation details and experiment results are included in the supplementary material.

**Baselines.** Here we design three baselines for extensive evaluations. The first is CIPS-3D [75], which naively fine-tunes the super-resolution  $G_1$  for view-consistent toonification. The second is E3DGE [32], which fine-tunes both  $G_0$  and  $G_1$  independently for true-3D toonification. Another prominent method is to leverage directional CLIP

loss [14, 2, 27] for adapting a pre-trained style-based generator. Here we extend StyleGAN-NADA [14] to 3D GAN and employ image CLIP directional loss for the evaluation. The baseline models are trained on the same dataset as DEFORMTOON3D, with each fine-tuning process applied to a single style. In contrast, DEFORMTOON3D is capable of accommodating all styles within a single model.

### 4.1. Comparisons with Baselines

**Qualitative Results.** We show the qualitative comparisons against the baselines in Fig. 4. DEFORMTOON3D fully captures the characteristics of the target domain with consistent identity preservations. CIPS-3D only provides texture-wise toonification and ignores geometry deformation. E3DGE suffers from mode collapse and tends to lose identity. StyleGAN-NADA fails to capture the characteristics of the target domain. Our method produces high-quality toonification with consistent identity preservations.

Table 1: **Quantitative evaluation over 10 styles.** DEFORMTOON3D achieves the best identity preservation (IP) and FID.

	CIPS-3D	E3DGE	NADA	Ours
IP $\uparrow$	0.681	0.707	0.535	<b>0.781</b>
FID $\downarrow$	50.6	34.0	59.3	<b>27.6</b>

**Quantitative Results.** To evaluate the fidelity and quality of toonification, we compare their identity preservation(IP) [11] and FID respectively in Tab. 1. In terms of IP, DEFORMTOON3D outperforms all baseline methods across the 10 styles provided, which underscores the benefits of retaining the 3D generator. In terms of FID, CIPS-3D [75] only fine-tunes  $G_1$  and achieves worse performance compared to E3DGE [32], which fine-tunes the whole generator. StyleGAN-NADA achieves the worst FID performance, which we attribute to the challenges of directly adopting 2D CLIP-based supervision on 3D GANs. A detailed breakdown by individual styles is available in the supplementary material.

Table 2: **User preference study.**

	CIPS-3D	E3DGE	NADA	Ours
Shape	20.8%	17.4%	3.4%	<b>58.4%</b>
Appearance	16.1%	16.4%	3.4%	<b>64.1%</b>
Identity	23.1%	7.1	5.4	<b>64.4%</b>
Overall	21.8%	9.9%	2.7%	<b>65.6%</b>

**User Study.** For the user preference study, we collected 2400 votes to select the preferred rendering results in terms of shape stylization, appearance stylization, identity preservation, and overall performance. As shown in Tab. 2, the proposed method gives the most preferable results despite the fact that the baseline methods are trained on a single style.

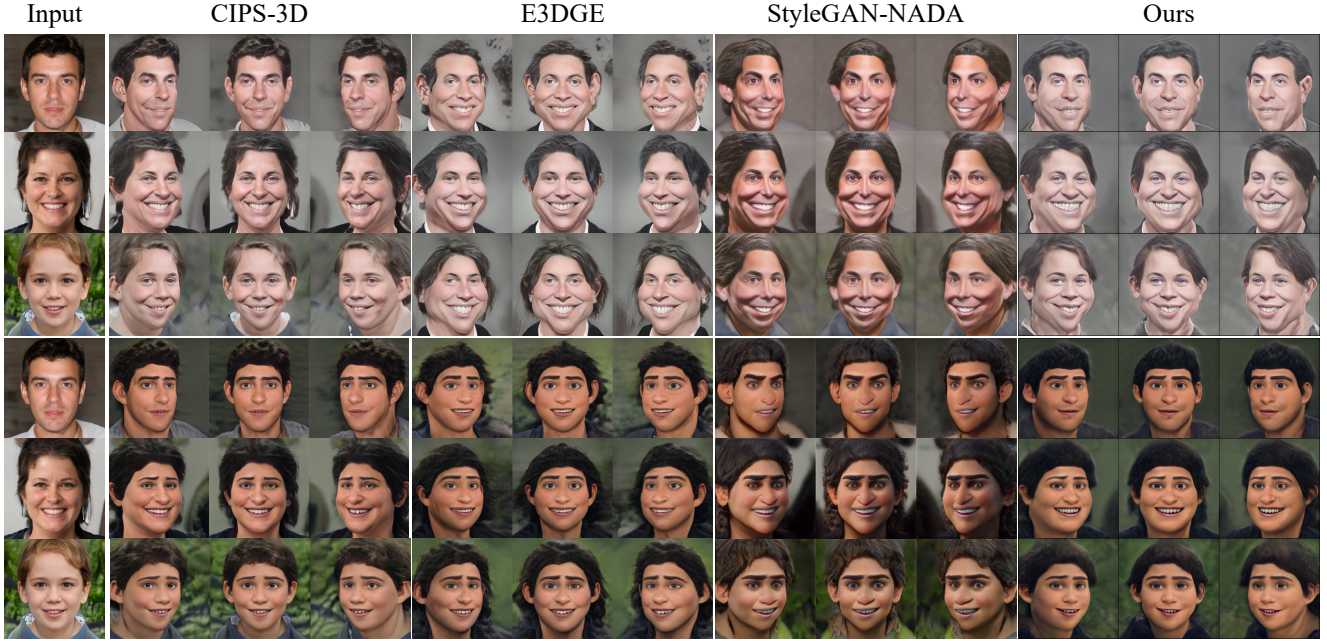


Figure 4: **Qualitative comparisons with baseline methods.** DEFORMTOON3D produces better performance against all baselines regarding toonification fidelity, diversity and identity preservation.

Table 3: **Storage cost comparison.** Values are averaged across 10 styles and shown in MB. DEFORMTOON3D achieves a considerably more efficient storage footprint against the baselines.

Methods	Trainable Params↓	Model Storage↓
CIPS-3D [75]	5.81	59.93
E3DGE [32]	7.64	76.4
StyleGAN-NADA [14]	7.64	76.4
Ours (single-style)	<b>3.82</b>	<b>11.46</b>

**Storage Cost Comparison.** We detail the storage costs in Table 3. Thanks to our unique design that disentangles geometry and texture, our method requires no fine-tuning and fewer parameters for training. In a single-style scenario, our method reduces storage needs by 85%. This advantage becomes even more pronounced with the multi-style version of  $H_D$ , achieving a storage saving of 98.5% in a 10-style scenario. This makes our approach particularly feasible for potential mobile applications.

## 4.2. Applications

To demonstrate the generality of full GAN latent space preservation, we show that DEFORMTOON3D could be easily extended to a series of downstream applications proposed for the original pre-trained GAN, including inversion, editing, and animation. To further validate DEFORMTOON3D’s unique geometry-texture decomposed toonification, we show results of flexible toonification style control. **Inversion and Editing.** With fully preserved 3D generative prior, DEFORMTOON3D could directly adopt pre-trained



Figure 5: **Editing of toonified results.** We show two attribute editing results over six styles. In row 1, we add bangs to the male identities, and in row 2, we add “Smile” to female identities. The edited results fully preserve the identity and abide by the style of the target domain.

3D GAN inversion framework and latent editing directions for semantic-aware editing over the toonified portrait. Here we adopt E3DGE [32] for 3D GAN inversion and show the inverted results of the real image in Fig. 6. With fully preserved GAN latent space, DEFORMTOON3D could be applied to real images over multiple styles with high quality. We also include the editing results of diverse styles in Fig. 5. As can be seen, our method produces consistent editing results and fully preserves the identity and the style of the target domain, regarding both the geometry and texture.

**Animatable Toonification.** Inspired by the success of 2D method [64] that obtains a rig-like control over StyleGAN-generated 2D faces, we propose a straightforward pipeline that aligns 3DMM [4] parameters with pre-trained 3D GAN latent space. Specifically, we train two MLP  $\mathcal{F}_w : \mathbb{R}^{|\mathcal{W}|} \rightarrow \mathbb{R}^{|\mathcal{M}|}$  and  $\mathcal{F}_m : \mathbb{R}^{|\mathcal{M}|} \rightarrow \mathbb{R}^{|\mathcal{W}|}$  to learn the bi-directional mapping between 3DMM parameter space  $\mathcal{M}$  and 3D GAN



Figure 6: **DEFORMTOON3D on the real images.** Our method enables multiple styles toonification with a single model, where both the texture and the geometry matches the target domain.

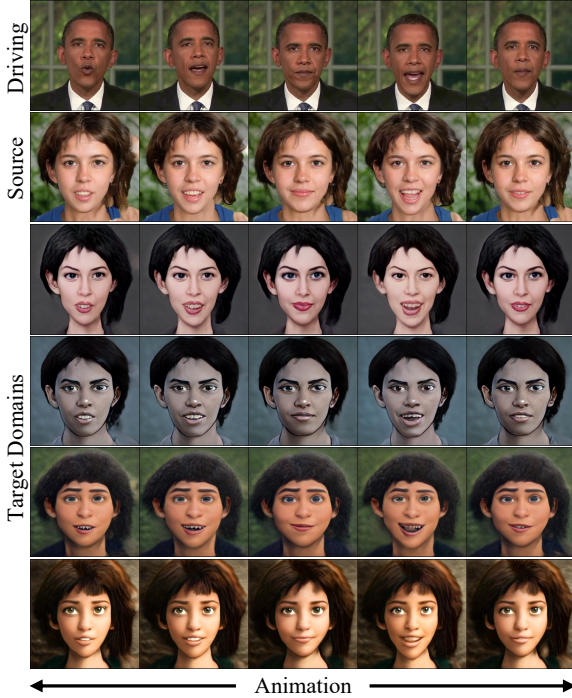


Figure 7: **Animation of stylized results.** We drive the real space images (row 2) with an input video (row 1). Since DEFORMTOON3D fully preserves the pre-trained GAN latent space, the driving direction of the real space could be directly applied to the style space (row 3 – 6). The expression on of animated toonified identities fully abide with the driving frames without affecting the toonification performance of target domain.

$\mathcal{W}$  space. To impose 3DMM-based control, given a latent code  $\mathbf{w}$ , we first infer its 3DMM parameters  $\tilde{\mathbf{m}} = \mathcal{F}_{\mathbf{w}}(\mathbf{w})$  and reconstructed 3DMM code  $\tilde{\mathbf{w}} = \mathcal{F}_{\mathbf{m}}(\tilde{\mathbf{m}})$ . After imposing 3DMM-based editing  $\tilde{\mathbf{m}}_{\Delta} = \tilde{\mathbf{m}} + \mathbf{m}_{\Delta}$ , we infer the corresponding edited  $\mathcal{W}$  code  $\tilde{\mathbf{w}}_{\Delta} = \mathcal{F}_{\mathbf{m}}(\tilde{\mathbf{m}}_{\Delta})$ . The final result is synthesized from  $\hat{\mathbf{w}} = \mathbf{w} + (\tilde{\mathbf{w}}_{\Delta} - \tilde{\mathbf{w}})$ . The whole framework is trained in a self-supervised manner with cycle consistency regularizations. Please refer to the supplementary material for more technical details and animation results.

We show the animated toonification results in Fig. 7. Here, we extract the 3DMM parameters from a driving video [18] using a pre-trained predictor [12]. The expression dimension of extracted parameters are injected to the sampled  $\mathbf{w}_R$  using the procedure described above. Four styles of animated toonified results are included. As can be seen, our designed animation pipeline could accurately drive the identity both in the real space (row 2) and the style spaces (row 3 – 6). The reenacted expressions are natural and abide with the driving input, which validates the generality of our method.

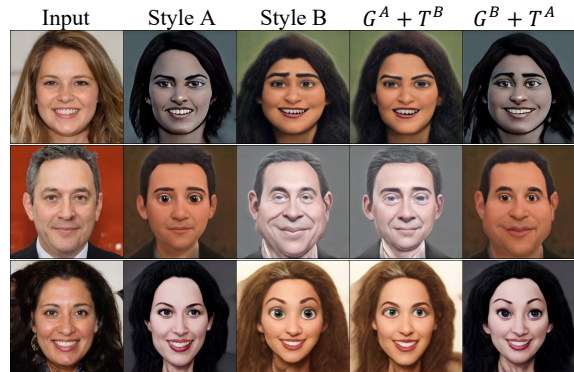


Figure 8: **Style swap.** Besides style interpolation, DEFORMTOON3D supports style swap of geometry and texture only across two styles. As shown here, given the real space input (col 1) and the toonification results of two spaces (col 2 – 3), DEFORMTOON3D could swap the geometry  $G$  and texture  $T$  of two styles independently (col 4 – 5), which cannot be achieved by previous methods.

**Toonification Style Control.** Due to the unique geometry-texture decomposition design, DEFORMTOON3D offers flexible style control. **First**, we achieve style degree control and show the results in Fig. 9. Geometry-wise, since  $H_D$  outputs the 3D deformation offsets  $\Delta \mathbf{x}_S$ , we simply interpolate the offsets with  $\tau * \Delta \mathbf{x}_S$  where  $\tau = 0$  represents an identical mapping of the real space. Texture-wise, we rescale the style mixing weight  $w$  of  $G_1$ , where  $w = 0$  preserves the color of the real space images. **Second**, due to the geometry-texture disentangled property, DEFORM-

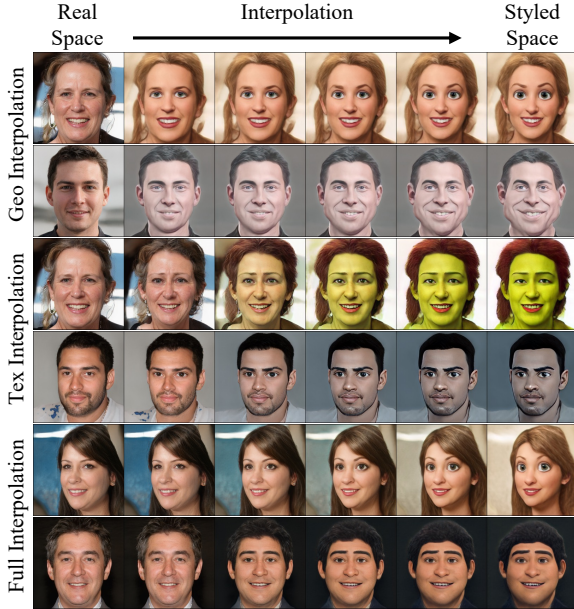


Figure 9: **Style degree control.** Thanks to the geometry-texture toonification decomposition, DEFORMTOON3D can specify geometry, texture, and full style control.

Table 4: Ablation study.

Method	IP $\uparrow$	FID $\downarrow$
w/o $\mathbf{w}_R$ condition	0.735	33.4
MLP as $H_D$	0.746	31.8
w/o $\mathcal{L}_{Adv}$	0.769	29.9
Ours	<b>0.781</b>	<b>27.6</b>

TOON3D naturally supports the toonification of shape and texture only. As shown in Fig. 8, we achieve the geometry-texture swap between multiple styles, where the geometry of one style could be combined with the texture of another style. This could not be achieved by all previous methods and opens up broader potential downstream applications.

### 4.3. Ablation Study

We ablate the effectiveness of our design choices in Tab. 4. **1) Instance code condition:** By removing instance code  $\mathbf{w}_R$ , StyleField tends to learn similar offsets for different identities, whereas adopting  $\mathbf{w}_R$  as the instance deformation conditions facilitates better toonification. **2) Style-Field architecture:** We replace the SIREN architecture with an MLP network as defined by D-NeRF [51] and observe the significant performance drop. This demonstrates the representation power of SIREN deformation field. **3) Adversarial training:** We also validate the effectiveness of adversarial loss, which brings noticeable improvement regarding both FID and identity similarity, respectively.

### 4.4. Limitations

As shown in Fig. 10 pertaining to the style ‘‘Comic’’ and ‘‘Slam Dunk’’, though DeformToon3D produces reasonable texture-wise stylization, the corresponding geometry still has noticeable artifacts. The proposed StyleField implicitly learns the correspondence between the paired data from the real space and the style space. Such correspondence is easier to learn with information cues such as illumination from the 3D-ish styles like Pixar or noticeable keypoints from caricature styles, but harder for styles with limited information cues like Comic.

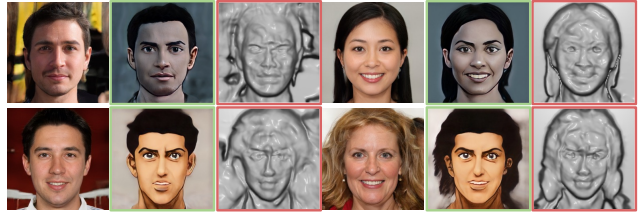


Figure 10: **Failure cases.**

## 5. Conclusion and Future Work

In this paper, we propose a novel 3D toonification framework DEFORMTOON3D for a fine-tuning free, geometry-texture decomposed 3D face toonification. We fully exploit the hierarchical characteristics of 3D GAN and introduce a StyleField to handle 3D geometry toonification of  $G_0$ , along with adaptive style mixing that injects texture information into  $G_1$ . Our method achieves high-quality toonification on both geometry and texture, outperforming existing methods. Thanks to the preservation of the 3D generative prior, DEFORMTOON3D facilitates a range of downstream applications.

As a pioneering effort in this field, we believe this work will inspire future works on free 3D toonification. First, to mitigate the geometry-texture ambiguity present in certain styles, introducing re-lighting during training could serve as a potential solution [43]. Second, a more flexible training paradigm could directly guide the 3D toonification process with a pre-trained vision-language model [53]. Third, future research could focus on integrating a comprehensive 3D animation pipeline [3, 63] into the toonification process. Moreover, the potential applicability of DeformToon3D to other 3D GANs [8] and shapes beyond human faces, such as full-body settings, is worth investigating.

**Acknowledgement.** This study is supported under the RIE2020 Industry Alignment Fund Industry Collaboration Projects (IAF-ICP) Funding Initiative, as well as cash and in-kind contribution from the industry partner(s). It is also supported by Singapore MOE AcRF Tier 2 (MOE-T2EP20221-0011, MOE-T2EP20221-0012) and NTU NAP Grant.



## A. Background

Since recent 3D-aware image generative models are all based on neural implicit representations, especially NeRF [37], here we briefly introduce the NeRF-based 3D representation and more StyleSDF details for clarification.

**NeRF-based 3D Representation.** NeRF [37] proposed an implicit 3D representation for novel view synthesis. Specifically, NeRF defines a scene as  $\{c, \sigma\} = F_{\Phi}(\mathbf{x}, \mathbf{v})$ , where  $\mathbf{x}$  is the query point,  $\mathbf{v}$  is the viewing direction from camera origin to  $\mathbf{x}$ ,  $c$  is the emitted radiance (RGB value),  $\sigma$  is the volume density. To query the RGB value  $C(\mathbf{r})$  of a point on a ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{v}$  shoot from the 3D coordinate origin  $\mathbf{o}$ , we have the volume rendering formulation,

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))c(\mathbf{r}(t), \mathbf{v})dt, \quad (5)$$

where  $T(t) = \exp(-\int_{t_n}^t \sigma(\mathbf{r}(s))ds)$  is the accumulated transmittance along the ray  $\mathbf{r}$  from  $t_n$  to  $t$ .  $t_n$  and  $t_f$  denote the near and far bounds.

**Hybrid 3D Generation.** In hybrid 3D generation [41, 8, 17], the intermediate feature map is calculated by replacing the color  $c$  with feature  $\mathbf{f}$ , namely  $\mathbf{F}(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{f}(\mathbf{r}(t), \mathbf{v})dt$ . Then, a StyleGAN [24, 25]-based decoder upsamples  $\mathbf{F}$  into high-resolution images with high-frequency details.

**SDF and Radiance-based Geometry Representation.** The intermediate geometry representation of  $G_0$  diversifies the characteristics of different 3D GANs. Specifically, StyleSDF [41] uses  $G_0$  to predict the signed distance  $d(\mathbf{x}) = G_0(\mathbf{w}, \mathbf{x})$  between the query point  $\mathbf{x}$  and the shape surface, where the density function  $\sigma(\mathbf{x})$  can be transformed from  $d(\mathbf{x})$  [41, 66, 71] for volume rendering [37]. The incorporation of SDF leads to higher-quality geometry in terms of expressiveness view consistency and clear definition of the surface.

In this paper, we mainly adopt StyleSDF [41] due to its high-quality geometry surface and high-fidelity texture. In StyleSDF, the Sigmoid activation function  $\sigma$  is replaced by  $\sigma(\mathbf{x}) = K_{\alpha}(d(\mathbf{x})) = \text{Sigmoid}(-d(\mathbf{x})/\alpha)/\alpha$ , where  $\alpha$  is a learned parameter that controls the tightness of the density around the surface boundary.

## B. Implementation Details

**CIPS-3D Baseline.** Following CLIPS-3D [75], we fine-tune  $G_1$  of StyleSDF on the toonified images with identical optimization parameters in the official implementation. The fine-tuning time for one style costs 10 V100 minutes.

**E3DGE Baseline.** Following E3DGE [32], we first fine-tune  $G_0$  for 400 iterations with batch size 24, and further fine-tune  $G_1$  for 400 iterations with batch size 8. All hyper-parameters are left unchanged with the official

StyleSDF [41] implementation. The overall fine-tuning time for one style costs around 30 minutes on a single V100 GPU.

**StyleGAN-NADA Baseline.** We reproduce StyleGAN-NADA [14] on StyleSDF with the following modifications. For  $G_0$  optimization, we fix the pre-trained mapping network, affine code transformations, view-direction MLP, color-prediction MLP, and density-prediction MLP. For  $G_1$  optimization, we follow the original implementation and fine-tune all weights except to RGB layers, affine code transformations, and mapping network. The  $k$  layers to optimize are also selected adaptively using StyleCLIP global loss. Other hyper-parameters and training procedures are left unchanged. The whole optimization costs around 5 minutes on a single V100 GPU.

### B.1. Additional Method Details

**StyleField.** Given the instance code  $\mathbf{w}$ , style code  $\mathbf{z}_S$ , we concatenate them along the channel dimension and send them into a 4-layer mapping network [24]. The mapping network first maps  $\mathbf{w} \oplus \mathbf{z}_S$  to a set of modulation signals  $\{\beta, \gamma\}$ , where  $\beta = \{\beta_i\}, \gamma = \{\gamma_i\}$ . To associate the given codes to the corresponding deformation, the modulation signals will be injected into the MLP network, serving as FiLM conditions [48, 13, 59] to modulate its features at different layers as  $\mathbf{f}_{i+1} = \sin(\gamma_i \cdot (\mathbf{W}_i \mathbf{f}_i + \mathbf{b}_i) + \beta_i)$ . To support multi-style code, we associate each style index with a learnable embedding. During inference, we pass in the corresponding style index to retrieve the style embedding for conditional deformation.

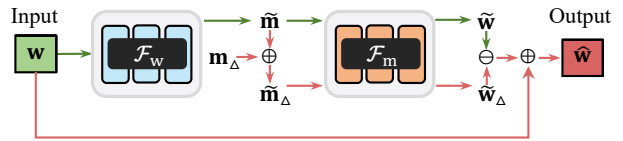


Figure 11: Inference pipeline of the proposed animation pipeline.

**Animatable Stylized Portrait Training Details.** We show the overall 3DMM animation inference pipeline in Fig. 11. Specifically, we train the whole framework in a self-supervised manner. In each iteration, we synthesize a batch of pose images  $\mathbf{I} = G(\mathbf{w})$ , where  $G = G_1 \circ G_0$ . For 3DMM supervision, we leverage the state-of-the-art 3DMM predictor [12] to infer the pseudo ground-truth 3DMM parameter  $\mathbf{m}_{GT}$ . With the synthesized training corpus, we reconstruct the input codes  $\hat{\mathbf{m}} = \mathcal{F}_w(\mathbf{w})$  and  $\hat{\mathbf{w}} = \mathcal{F}_m(\mathbf{m}_{GT})$  and impose MSE reconstruction loss. We further render the reconstructed code  $\hat{\mathbf{I}} = G(\hat{\mathbf{w}})$  and  $\hat{\mathbf{I}}_{\mathcal{M}} = \text{DFR}(\hat{\mathbf{m}})$ , where DFR is a differentiable render [52] that renders the reconstructed 3DMM mesh to image. The rendered images are supervised with corresponding loss [12], which yields bet-

ter performance in our observations.

This training objective shall guarantee plausible 3DMM editing using the procedure stated in the main context. However, in practice, we find the training is unstable and predicted codes are not disentangled well during inference. We make the following modifications to the overall training pipeline and improve the editing performance:

First, we observe that the 3DMM head pose parameters, including a head rotation  $\mathbf{R} \in SO(3)$  and translation  $\mathbf{t} \in \mathbb{R}^3$  parameters, are in contradict with the pose control of 3D GAN. This deteriorates the disentanglement of learned codes and destabilizes the training since the predicted 3DMM codes  $\tilde{\mathbf{m}}$  must contain accurate head pose to minimize the reconstruction loss with  $\text{DFR}(\mathbf{m})$ . To address this issue, we mask out the head rotation  $\mathbf{R}$  and translation  $\mathbf{t} \in \mathbb{R}^3$  dimension in all 3DMM parameters with the binary mask. This enforces all the 3DMM images  $\mathbf{I}_{\mathcal{M}}$  to be rendered from frontal pose and encourages the networks to focus on facial expression  $\delta \in \mathbb{R}^{64}$  alignment between  $\mathcal{W}$  and  $\mathcal{M}$ .

Second, to further impose identity preservation and bijective mapping between two spaces, we introduce cycle training [76] which regularizes  $\mathbf{w} \approx \mathcal{F}_{\mathbf{w}}(\tilde{\mathbf{m}})$  and  $\mathbf{m} \approx \mathcal{F}_{\mathbf{m}}(\tilde{\mathbf{w}})$ . The cycle loss is also imposed on the image space.

Third, to imitate the inference pipeline, in each training iteration, we randomly shuffle the expression dimension of all the 3DMM code  $\mathbf{m}_{\text{GT}}$  within a batch and generate a new set of codes  $\tilde{\mathbf{m}}_{\text{GT}}$ . The rendered image from  $\mathcal{F}_{\mathbf{m}}(\tilde{\mathbf{m}}_{\text{GT}})$  shall maintain the same identity with  $\mathbf{I}$  with identical pose of  $\text{DFR}(\tilde{\mathbf{m}}_{\text{GT}})$ . We impose the identity preservation loss [11] and landmark loss over the rendered 3DMM image [12] as supervisions. This strategy further reduces the domain gap between training and inference and further improves the final editing performance.

Fourth, we further leverage the style-based hierarchical structure within StyleSDF and reduces the attribute entanglement. Specifically, rather than using the edited  $\mathcal{W}$  code  $\tilde{\mathbf{w}}_{\Delta}$  for all the style layers in  $G$ , we conduct layer-wise editing effect analysis and find that only the first 2 layers of  $G_0$  will handle the expression-relevant information of the synthesized image. Using the edited code for later layers will result in other attributes editing, *e.g.*, adding glasses or changing the hair structure. Therefore, we leave the remaining 7 layers of  $G_0$  and all 10 layers in  $G_1$  unchanged and only use the edited code for the top 2  $G_0$  layers. This yields better disentanglement during the 3DMM-controlled style editing.

Training-wise, we adopt identical MLP architecture from PixelNeRF [72] to implement both  $\mathcal{F}_{*}$  networks and adopt a batch size of 4 with learning rate  $5 \times 10^{-4}$  during the optimization. The networks are trained for 50,000 iterations, which costs around 2 days on a single V100 GPU. Please refer to the released code for more details.

## B.2. Additional Ablation Study

The robustness of the number of style codes is ablated in Table 5. Experiments are conducted with 1, 2, and 5 styles per model so that under each setting the 10 styles can be evenly divided into different runs for the sake of comparison.

Table 5: Ablation on the number of styles.

# styles per model	1	2	5	10
Identity similarity $\uparrow$	0.795	0.776	0.784	0.781
FID $\downarrow$	27.5	28.1	27.9	27.6

The effectiveness of the elastic loss is also ablated qualitatively. As shown in the visualized mesh in Fig. 12, the elastic loss is effective in preventing discontinuous deformation and leads to smoother geometry in the styled space.

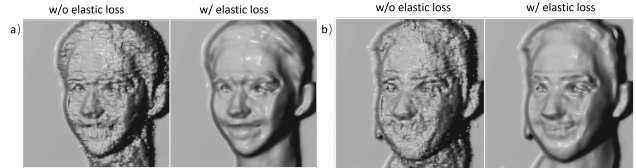


Figure 12: Ablation on elastic loss. Without *v.s.* with elastic loss for a) female and b) male, respectively.

## B.3. Additional Quantitative and Qualitative Results

The detailed breakdown of toonification fidelity and quality are shown in Tab. 6 and Tab. 7 respectively. We include more qualitative experiment results here. In Fig. 13 we include more comparisons with the baseline methods,

Table 6: Quantitative evaluation in terms of identity similarity $\uparrow$ . DEFORMTOON3D achieves the best identity consistency over all the 10 styles.

Domains	CIPS-3D	E3DGE	NADA	Ours
Pixar	0.765	0.748	0.564	<b>0.812</b>
Comic	0.643	0.614	0.496	<b>0.729</b>
Slam Dunk	0.672	0.765	0.552	<b>0.780</b>
Caricature I	0.648	0.592	0.455	<b>0.708</b>
Caricature II	0.655	0.698	0.538	<b>0.785</b>
Caricature III	0.637	0.644	0.495	<b>0.725</b>
Croods	0.796	0.831	0.626	<b>0.860</b>
Shrek	0.708	0.794	0.599	<b>0.835</b>
Rapunzel	0.603	0.696	0.564	<b>0.782</b>
Hiccup	0.684	0.688	0.464	<b>0.796</b>
Average	0.681	0.707	0.535	<b>0.781</b>

Table 7: **Quantitative evaluation in terms of FID↓**. DEFORMTOON3D achieves the best FID over 9 of the 10 styles.

Domains	CIPS-3D	E3DGE	NADA	Ours
Pixar	33.6	36.8	39.9	<b>21.5</b>
Comic	61.9	44.8	70.9	<b>33.3</b>
Slam Dunk	78.1	41.8	75.9	<b>37.3</b>
Caricature I	28.7	30.1	52.9	<b>16.0</b>
Caricature II	76.7	58.1	102.6	<b>56.4</b>
Caricature III	47.1	<b>25.8</b>	54.8	27.2
Croods	36.9	30.9	58.5	<b>22.5</b>
Shrek	36.0	32.1	47.2	<b>20.3</b>
Rapunzel	65.0	30.5	44.1	<b>17.2</b>
Hiccup	42.3	32.8	46.6	<b>24.6</b>
Average	50.6	34.0	59.3	<b>27.6</b>

which demonstrates that DEFORMTOON3D produces better quality against existing methods. In Fig. 14 we show more toonification results over real images. The proposed methods yield plausible results with consistent identity preservations. We further include the stylized texture and shape pair in Fig. 15 and validate that our method produces high-quality stylization over both texture and shape.



Figure 13: **Additional qualitative comparisons with baseline methods.** DEFORMTOON3D produces better performance against all baselines regarding toonification fidelity, diversity, and identity preservation. Better zoom in.



Figure 14: **Additional results of DEFORMTOON3D on the real images.** Our method enables multiple styles toonification with a single model, where both the texture and the geometry matches the target domain.

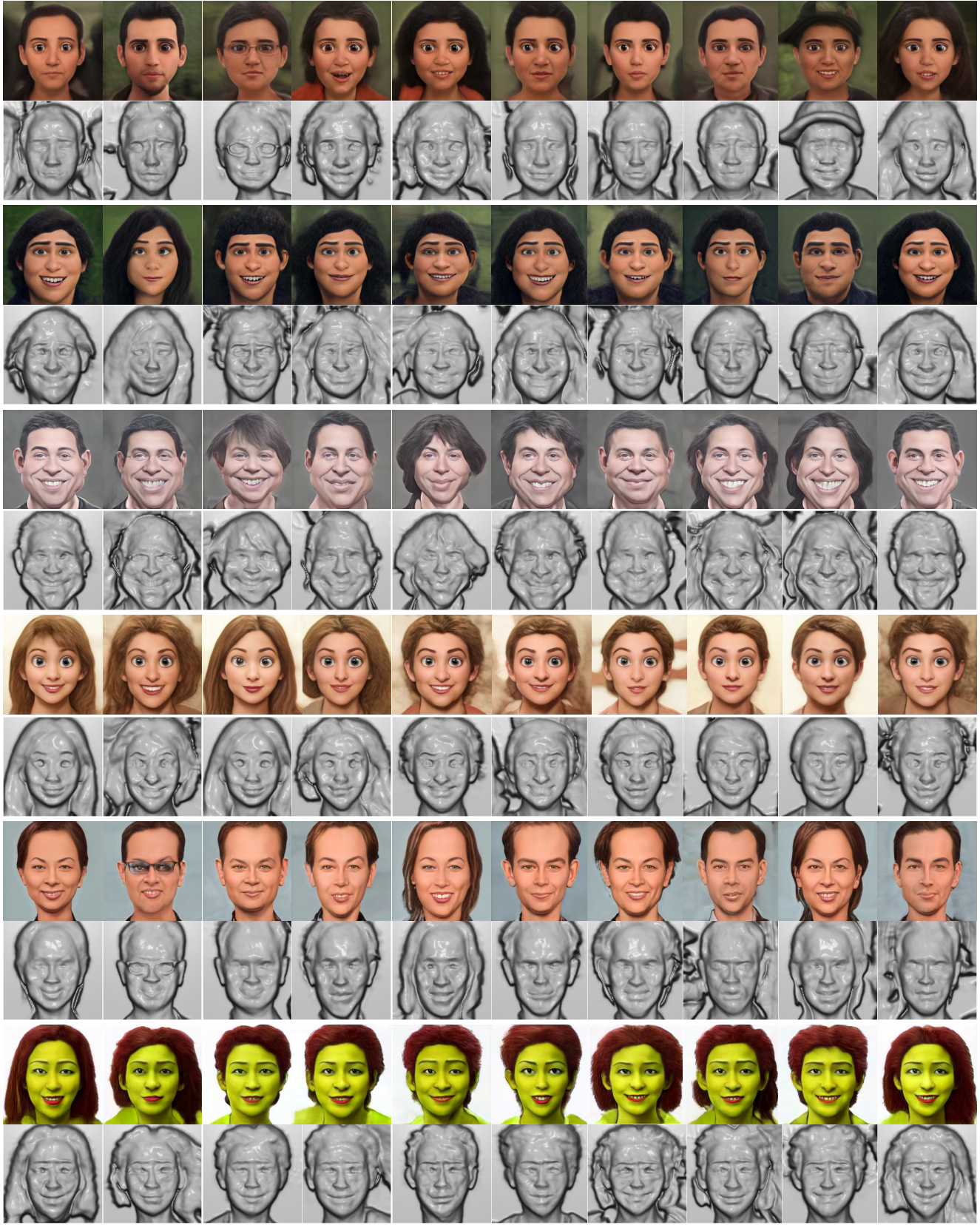


Figure 15: Additional results of DEFORMTOON3D with stylized texture and shape.

## References

- [1] Rameen Abdal, Hsin-Ying Lee, Peihao Zhu, Menglei Chai, Aliaksandr Siarohin, Peter Wonka, and S. Tulyakov. 3DAvatarGAN: Bridging domains for personalized editable avatars. *arXiv*, abs/2301.02700, 2023. 2, 3
- [2] Aibek Alanov, Vadim Titov, and Dmitry Vetrov. HyperDomainNet: Universal domain adaptation for generative adversarial networks. 2022. 5
- [3] Alexander W. Bergman, Petr Kellnhofer, Yifan Wang, Eric Chan, David B. Lindell, and Gordon Wetzstein. Generative neural articulated radiance fields. In *NIPS*, 2022. 8
- [4] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3D faces. In *SIGGRAPH*, 1999. 6
- [5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*. OpenReview.net, 2019. 2
- [6] Kaidi Cao, Jing Liao, and Lu Yuan. CariGANs: unpaired photo-to-caricature translation. *ACM TOG*, 37(6):1–14, 2018. 2
- [7] Eric Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and G. Wetzstein. Pi-GAN: Periodic Implicit Generative Adversarial Networks for 3D-Aware Image Synthesis. In *CVPR*, 2021. 2
- [8] 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 *CVPR*, 2022. 2, 3, 8, 9
- [9] Wenzheng Chen, Jun Gao, Huan Ling, Edward J. Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3D objects with an interpolation-based differentiable renderer. 2
- [10] Min Jin Chong and David Forsyth. GANs N’ Roses: Stable, controllable, diverse image to image translation. *arXiv preprint arXiv:2106.06561*, 2021. 2
- [11] Jiankang Deng, Jia Guo, Xue Niannan, and Stefanos Zafeiriou. ArcFace: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 5, 10
- [12] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3D face reconstruction with weakly-supervised learning: From single image to image set. In *CVPR*, 2019. 7, 9, 10
- [13] Vincent Dumoulin, Ethan Perez, Nathan Schucher, Florian Strub, Harm de Vries, Aaron Courville, and Yoshua Bengio. Feature-wise transformations. *Distill*, 2018. <https://distill.pub/2018/feature-wise-transformations>. 9
- [14] Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, and Daniel Cohen-Or. StyleGAN-NADA: Clip-guided domain adaptation of image generators. *arXiv*, abs/2108.00946, 2021. 3, 5, 6, 9
- [15] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR*, pages 2414–2423, 2016. 2
- [16] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014. 2
- [17] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. StyleNeRF: A style-based 3D-aware generator for high-resolution image synthesis. In *ICLR*, 2021. 3, 9
- [18] Yudong Guo, Keyu Chen, Sen Liang, Yongjin Liu, Hujun Bao, and Juyong Zhang. AD-NeRF: Audio driven neural radiance fields for talking head synthesis. In *ICCV*, 2021. 7
- [19] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Escaping plato’s cave: 3D shape from adversarial rendering. In *ICCV*, 2019. 2
- [20] Fangzhou Hong, Zhaoxi Chen, Yushi Lan, Liang Pan, and Ziwei Liu. EVA3D: Compositional 3D human generation from 2D image collections. In *ICLR*, 2023. 3
- [21] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, pages 1501–1510. 4
- [22] Wonjoon Jin, Nuri Ryu, Geon-Yeong Kim, Seung-Hwan Baek, and Sunghyun Cho. Dr3D: Adapting 3D GANs to artistic drawings. *SIGGRAPH Asia 2022 Conference Papers*, 2022. 2, 3
- [23] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *ICLR*, 2018. 3, 5
- [24] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 2, 3, 4, 5, 9
- [25] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *CVPR*, 2020. 2, 3, 9
- [26] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *CVPR*, 2020. 2, 3, 4
- [27] Gwanghyun Kim and Se Young Chun. DATID-3D: Diversity-preserved domain adaptation using text-to-image diffusion for 3D generative model, 2022. 2, 3, 5
- [28] Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwang Hee Lee. U-GAT-IT: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. In *ICLR*, 2019. 2
- [29] Sunnie SY Kim, Nicholas Kolkin, Jason Salavon, and Gregory Shakhnarovich. Deformable style transfer. In *ECCV*, 2020. 2
- [30] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, volume abs/1412.6980, 2015. 5
- [31] Yushi Lan, Chen Change Loy, and Bo Dai. DDF: Correspondence Distillation from NeRF-based GAN. *arXiv preprint arXiv:2212.09735*, 2022. 3
- [32] Yushi Lan, Xuyi Meng, Shuai Yang, Chen Change Loy, and Bo Dai. E3DGE: Self-supervised geometry-aware encoder for style-based 3D gan inversion. In *CVPR*, 2022. 2, 3, 4, 5, 6, 9
- [33] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. MaskGAN: Towards diverse and interactive facial image manipulation. In *CVPR*, 2020. 3
- [34] Bing Li, Yuanlue Zhu, Yitong Wang, Chia-Wen Lin, Bernard Ghanem, and Linlin Shen. AniGAN: Style-guided generative adversarial networks for unsupervised anime face generation. *IEEE TMM*, 2021. 2

- [35] Jing Liao, Yuan Yao, Lu Yuan, Gang Hua, and Sing Bing Kang. Visual attribute transfer through deep image analogy. *ACM Transactions on Graphics*, 36(4):120, 2017. 2
- [36] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *CVPR*, June 2019. 2
- [37] 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. 2, 3, 9
- [38] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yongliang Yang. HoloGAN: Unsupervised Learning of 3D Representations From Natural Images. In *ICCV*, 2019. 2
- [39] Michael Niemeyer and Andreas Geiger. GIRAFFE: Representing scenes as compositional generative neural feature fields. In *CVPR*, 2021. 3
- [40] Utkarsh Ojha, Yijun Li, Jingwan Lu, Alexei A. Efros, Yong Jae Lee, Eli Shechtman, and Richard Zhang. Few-shot image generation via cross-domain correspondence. *CVPR*, pages 10738–10747, 2021. 3
- [41] 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 *CVPR*, 2021. 2, 3, 5, 9
- [42] Xingang Pan, Bo Dai, Ziwei Liu, Chen Change Loy, and Ping Luo. Do 2D GANs know 3D shape? Unsupervised 3D Shape Reconstruction from 2D Image GANs. In *ICLR*, 2021. 2
- [43] Xingang Pan, Xudong Xu, Chen Change Loy, Christian Theobalt, and Bo Dai. A Shading-Guided Generative Implicit Model for Shape-Accurate 3D-Aware Image Synthesis. In *NIPS*, 2021. 8
- [44] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *CVPR*. IEEE, 2019. 2
- [45] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, 2021. 3, 4
- [46] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *TOG*, 40(6), dec 2021. 3
- [47] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. StyleCLIP: Text-driven manipulation of stylegan imagery. *ICCV*, pages 2065–2074, 2021. 3
- [48] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. FiLM: Visual reasoning with a general conditioning layer. In *AAAI*, volume 32, 2018. 9
- [49] Justin NM Pinkney and Doron Adler. Resolution dependent gan interpolation for controllable image synthesis between domains. *arXiv preprint arXiv:2010.05334*, 2020. 2
- [50] Justin N. M. Pinkney and Doron Adler. Resolution dependent gan interpolation for controllable image synthesis between domains. *arXiv*, abs/2010.05334, 2020. 2, 3
- [51] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *CVPR*, 2020. 8
- [52] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 9
- [53] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *CVPR*, pages 10674–10685, 2021. 3, 8
- [54] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3D-aware image synthesis. In *NIPS*, 2020. 2
- [55] Ahmed Selim, Mohamed Elgharib, and Linda Doyle. Painting style transfer for head portraits using convolutional neural networks. *ACM TOG*, 35(4):1–18, 2016. 2
- [56] Yujun Shen, Ceyuan Yang, Xiaou Tang, and Bolei Zhou. InterFaceGAN: Interpreting the disentangled face representation learned by GANs. *PAMI*, PP, 2020. 2, 3
- [57] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. In *CVPR*, 2021. 3
- [58] Yichun Shi, Debayan Deb, and Anil K Jain. WarpGAN: Automatic caricature generation. In *CVPR*, pages 10762–10771, 2019. 2
- [59] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *NIPS*, 2020. 4, 9
- [60] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene Representation Networks: Continuous 3D-structure-aware neural scene representations. In *NIPS*, 2019. 2
- [61] Guoxian Song, Linjie Luo, Jing Liu, Wan-Chun Ma, Chun-Pong Lai, Chuanxia Zheng, and Tat-Jen Cham. AgileGAN: stylizing portraits by inversion-consistent transfer learning. *ACM Transactions on Graphics (TOG)*, 40:117:1–117:13, 2021. 3
- [62] Kunpeng Song, Ligong Han, Bingchen Liu, Dimitris Metaxas, and Ahmed Elgammal. Diffusion guided domain adaptation of image generators. *arXiv preprint https://arXiv.org/abs/2212.04473*, 2022. 3
- [63] Jingxiang Sun, Xuan Wang, Lizhen Wang, Xiaoyu Li, Yong Zhang, Hongwen Zhang, and Yebin Liu. Next3D: Generative neural texture rasterization for 3D-aware head avatars. In *CVPR*, 2023. 8
- [64] 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 *CVPR*, pages 6142–6151, 2020. 2, 6
- [65] Ayush Tewari, R. MallikarjunB., Xingang Pan, Ohad Fried, Maneesh Agrawala, and Christian Theobalt. Disentangled3D: Learning a 3d generative model with disentangled



- geometry and appearance from monocular images. *CVPR*, pages 1506–1515, 2022. 3
- [66] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. NeuS: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In *NeurIPS*, 2021. 9
- [67] Erroll Wood, Tadas Baltrušaitis, Charlie Hewitt, Sebastian Dziadzio, Matthew Johnson, Virginia Estellers, Thomas J. Cashman, and Jamie Shotton. Fake it till you make it: face analysis in the wild using synthetic data alone. *ICCV*, pages 3661–3671, 2021. 4
- [68] Zongze Wu, Yotam Nitzan, Eli Shechtman, and Dani Lischinski. StyleAlign: Analysis and applications of aligned stylegan models. *arXiv preprint arXiv:2110.11323*, 2021. 2
- [69] Shuai Yang, Liming Jiang, Ziwei Liu, , and Chen Change Loy. VToonify: Controllable high-resolution portrait video style transfer. *ACM Transactions on Graphics (TOG)*, 41(6):1–15, 2022. 3
- [70] Shuai Yang, Liming Jiang, Ziwei Liu, and Chen Change Loy. Pastiche master: Exemplar-based high-resolution portrait style transfer. In *CVPR*, 2022. 2, 3, 4, 5
- [71] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *NIPS*, 2021. 9
- [72] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. PixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 10
- [73] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 4
- [74] Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso, Antonio Torralba, and Sanja Fidler. DatasetGAN: Efficient labeled data factory with minimal human effort. In *CVPR*, 2021. 4
- [75] Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. CIPS-3D: A 3D-Aware Generator of GANs Based on Conditionally-Independent Pixel Synthesis. *arXiv*, 2021. 3, 5, 6, 9
- [76] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 10